

CSE 5523: Machine Learning Setup



THE OHIO STATE UNIVERSITY

Questions?

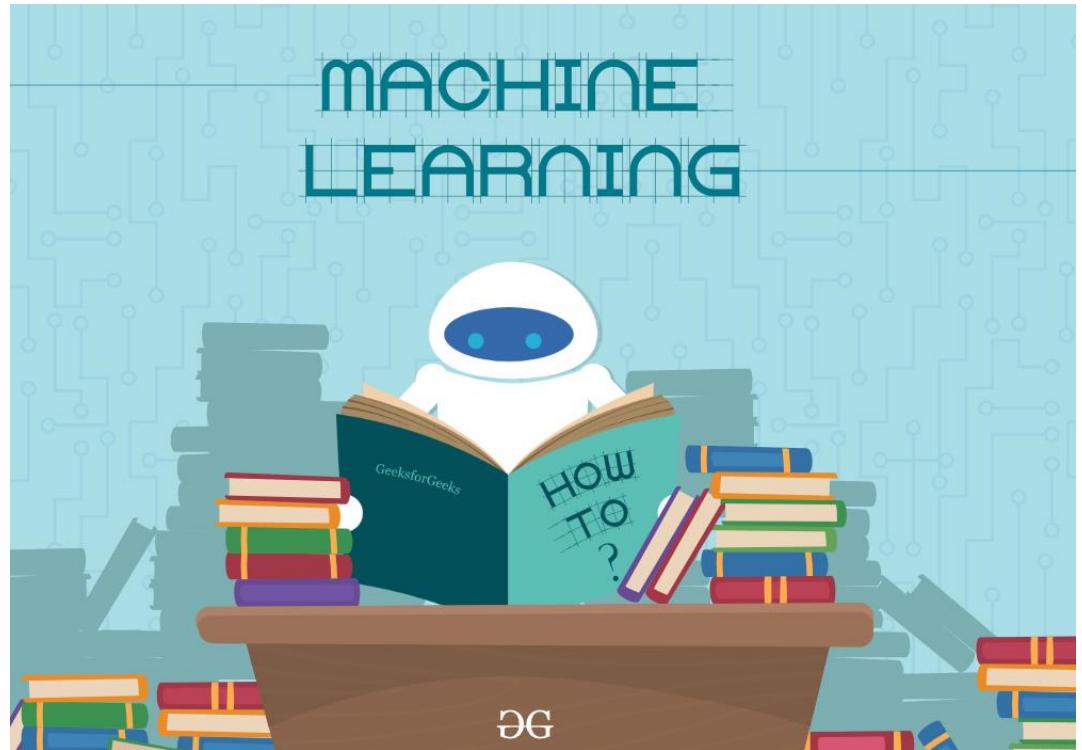


THE OHIO STATE UNIVERSITY

Today

Machine learning setup

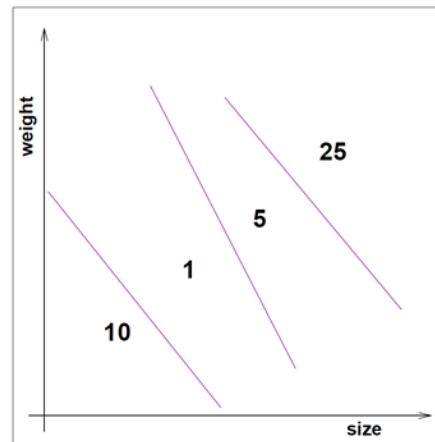
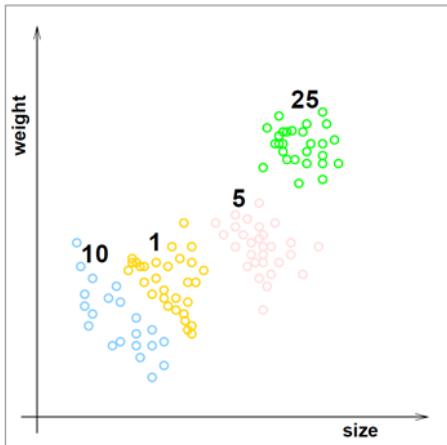
- Review
- Course overview
- More formal definitions



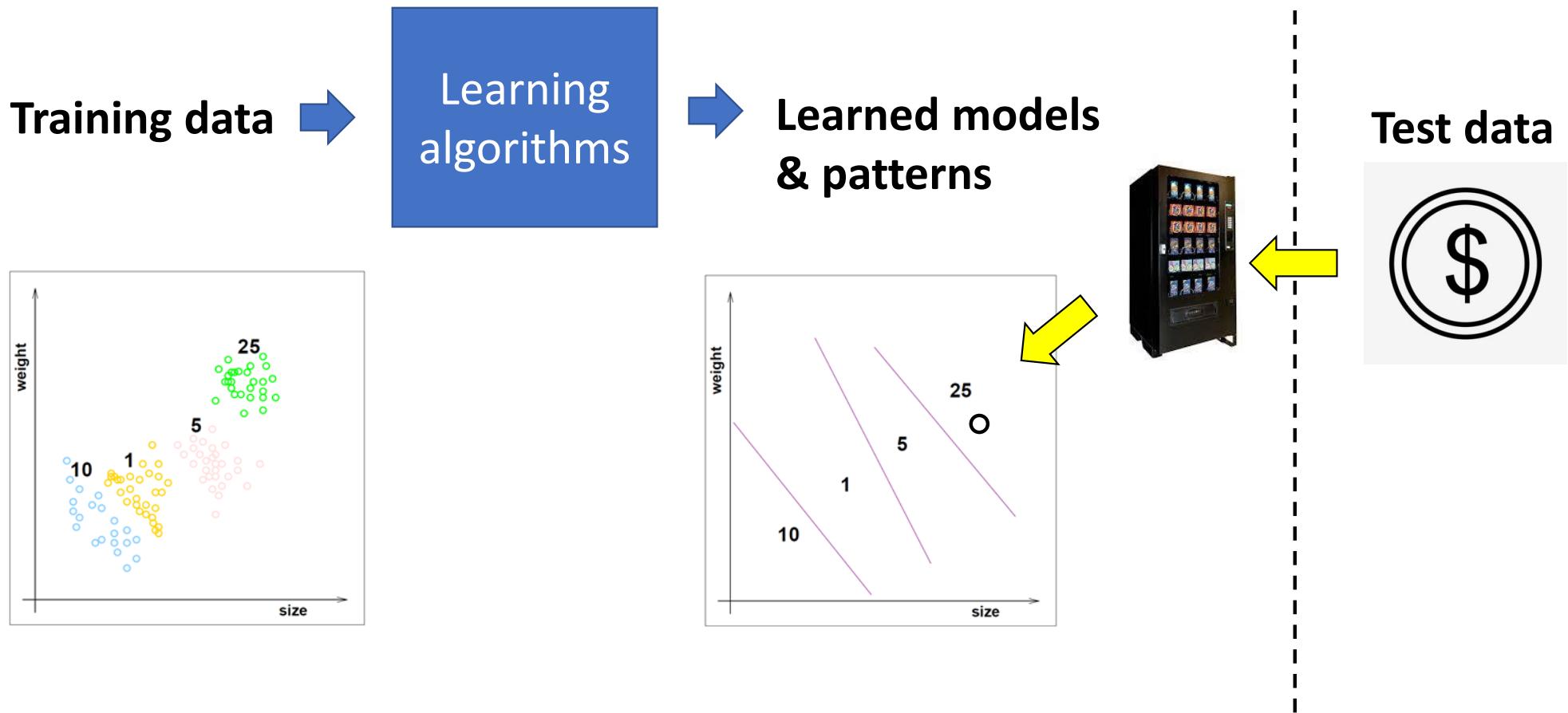
What is machine learning?



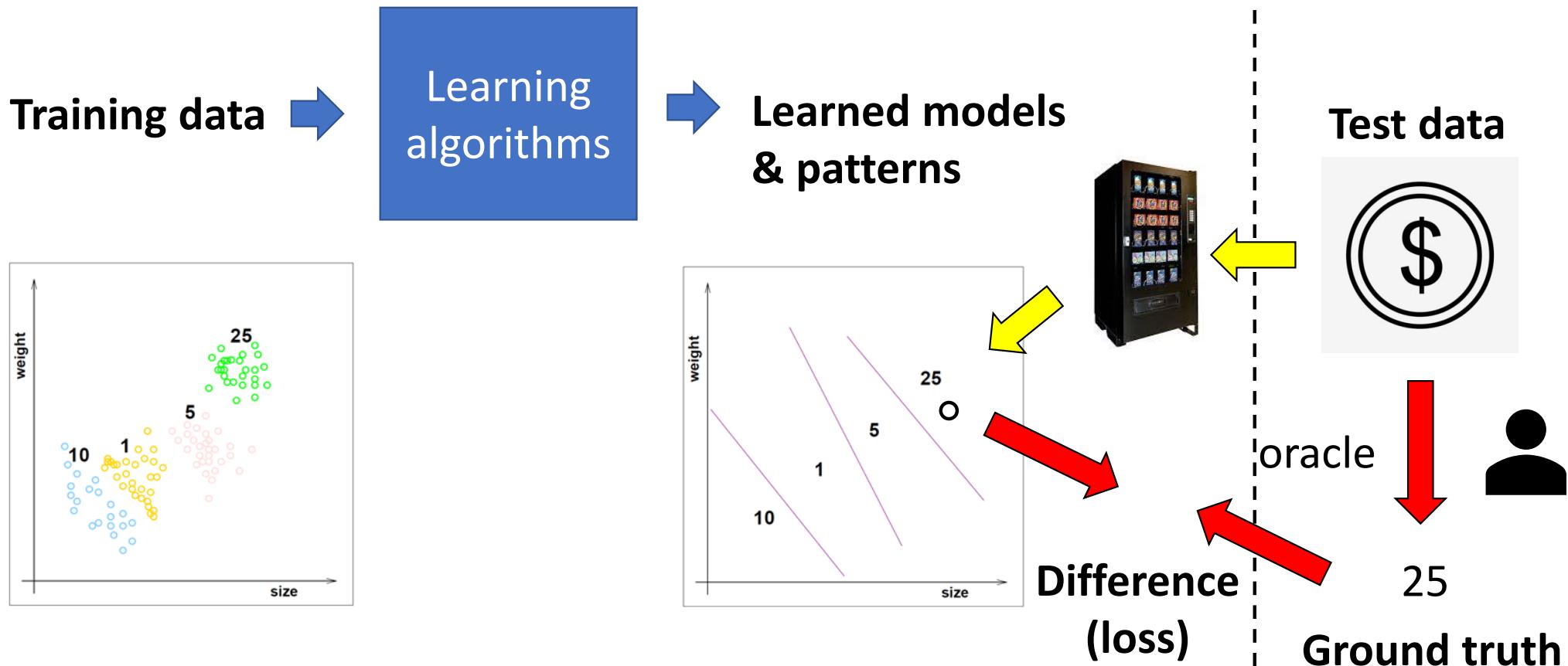
“Supervised” ML pipeline



“Supervised” ML pipeline

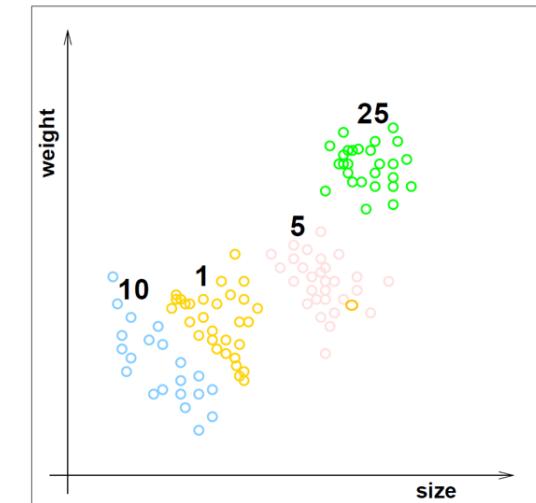


How to know if the model works well or not?



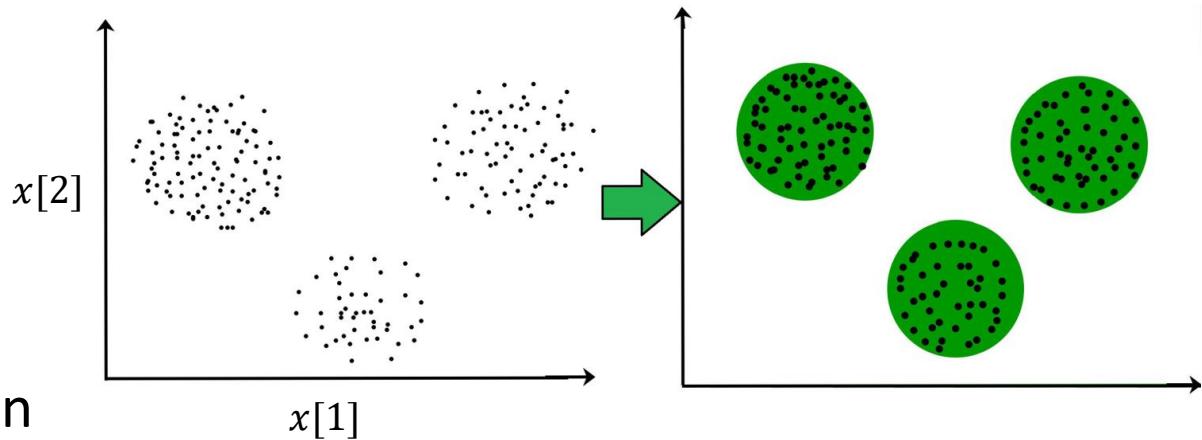
Different flavors of ML problems

- Several paradigms (e.g., dependent on the types of “*training*” data):
 - **Supervised learning (w/ labels)**
 - Training data: $\{(x_1, y_1), \dots, (x_N, y_N)\}$
 - x_n (or just x) is the input instance/example/feature vector $\in \mathbb{R}^d$ (or \mathcal{X} : feature space)
 - y_n (or just y) is the label, e.g., a real number $\in \mathbb{R}$ or a category index $\in \mathcal{C}$ (or \mathcal{Y} : label space)
 - Aim to predict (as in many of the previous list of applications)
 - **Unsupervised learning (w/o labels)**
 - Training data: $\{x_1, \dots, x_N\}$
 - Aim to discover hidden and latent patterns and explore data
 - Reinforcement learning
 - Many others

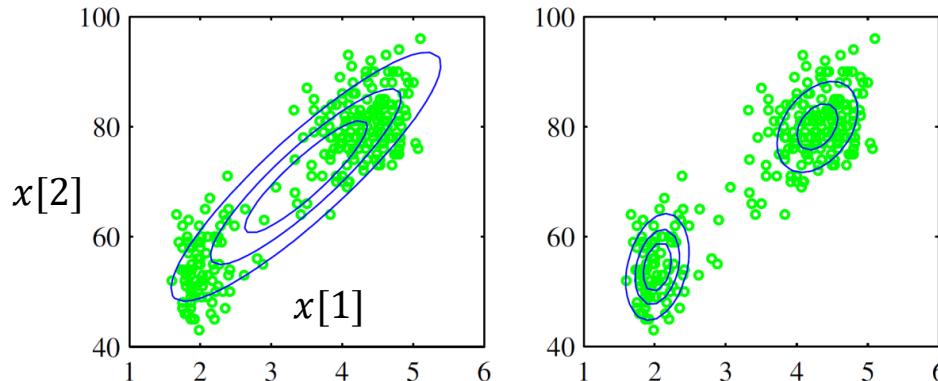


Unsupervised learning

- Clustering:



- Distribution/Density estimation
 - Can sample from the distribution to generate new data

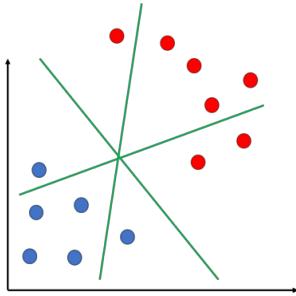


Classification vs. regression

- Classification
 - Form: Feature vector → Desired category (integer number index)
 - Ex: 1: Cat, 2: Dog, 3: Horse,, 1000: Car
- Regression (curve Fitting)
 - Form: Feature vector → Desired real number (e.g., price, chance, etc.)

Machine learning: capture patterns from training data
that can be generalized to future unseen data

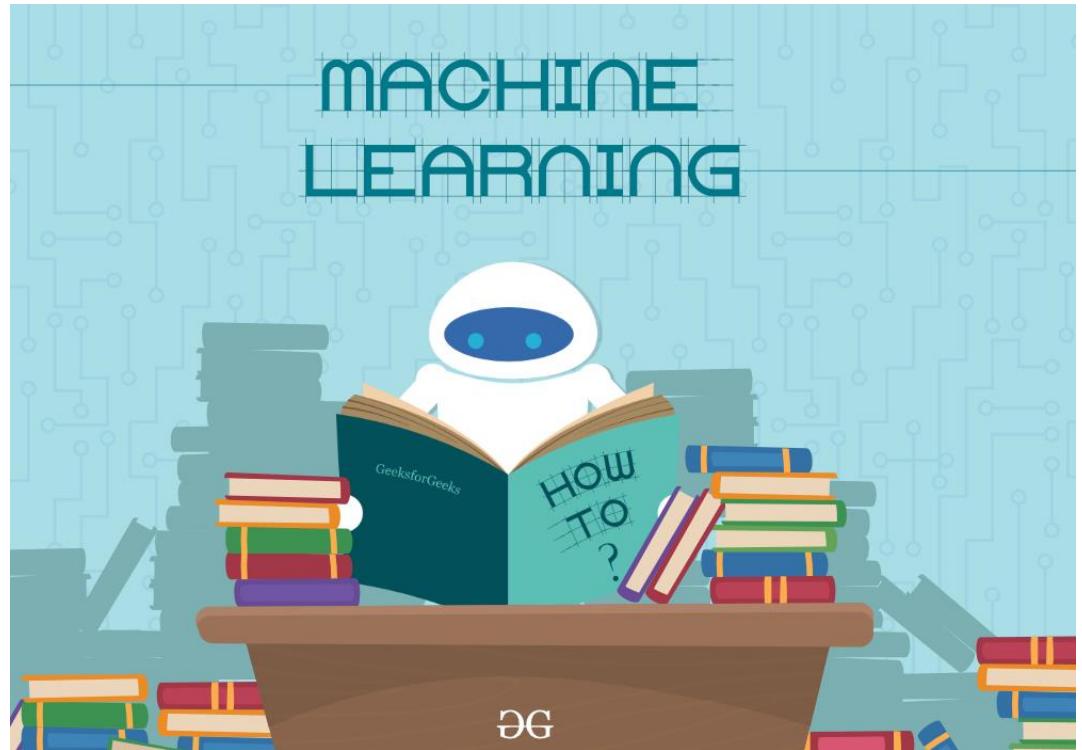
Questions to ML

- What assumptions do we need for learning to be possible?
 - Data and Models
- Given a fixed model type (e.g., linear boundaries), how can the machine output the “right” model?A 2D coordinate system with a horizontal x-axis and a vertical y-axis. There are seven blue dots clustered in the lower-left quadrant and six red dots clustered in the upper-right quadrant. Two parallel green lines are drawn through the data, representing different linear models. The top line separates all blue dots from all red dots. The bottom line separates most blue dots from most red dots, but it misclassifies one blue dot located above and to the right of the line.
- How many “training” data samples are needed to ensure that the output model will generalize well to the unseen test data?
- What is the practice of solving an ML problem?

Today

Machine learning setup

- Review
- Course overview
- More formal definitions



Topics

- **Supervised learning**
 - Regression
 - Classification (nearest neighbors, linear models, SVM)
 - Probabilistic approaches
 - Ensemble approaches
 - Kernel methods
 - Optimization
- **Machine learning foundation, theories, practice**
 - Empirical risk minimization and regularization
 - PAC learning and VC dimension
 - Bias and variance, over-fitting vs under-fitting
 - Debugging

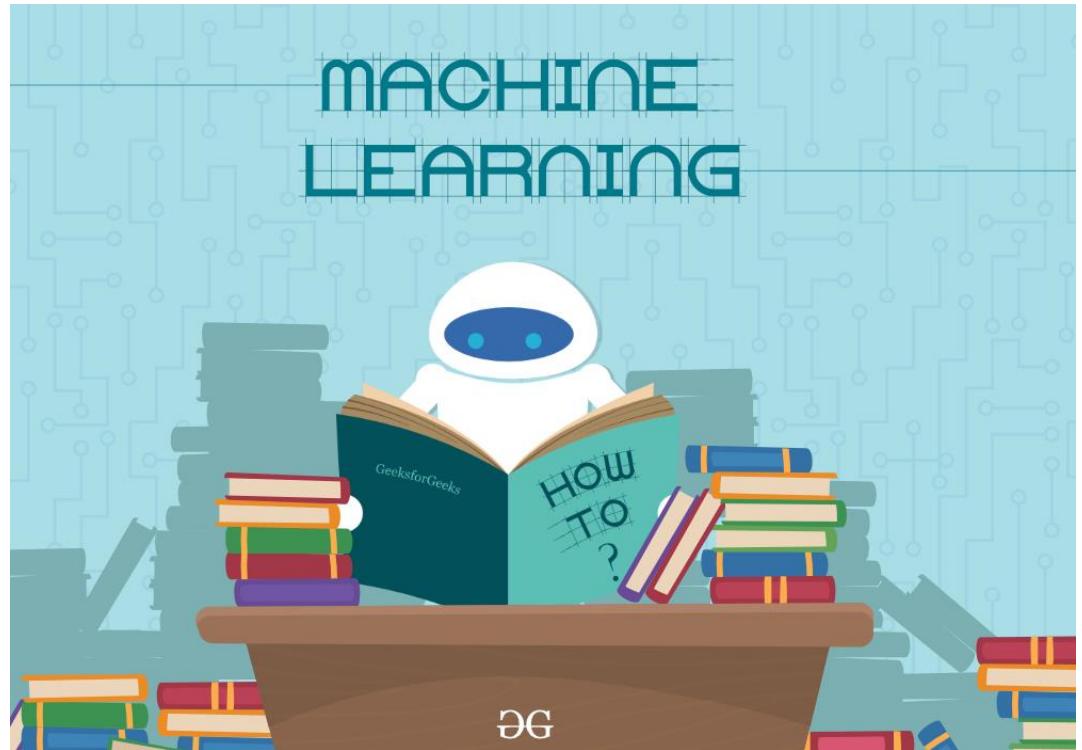
Topics

- **Unsupervised learning**
 - K-means clustering
 - GMM and EM
 - Dimensionality reduction and metric learning (supervised or unsupervised)
- **Neural networks and deep learning**
 - Backpropagation and optimization
 - MLP
 - CNN
 - Transformer
- **Optional: Advanced ML paradigms or approaches**
 - Generative models: VAE, GAN, diffusion models
 - Domain adaptation, semi-supervised learning, self-supervised learning

Today

Machine learning setup

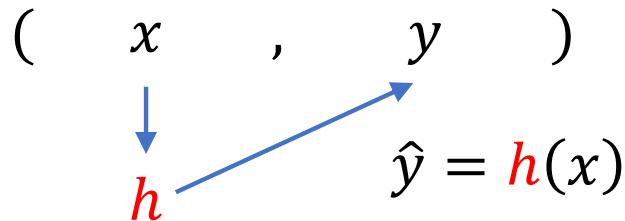
- Review
- Course overview
- More formal definitions



More formal definition of supervised learning

- An “unknown” **true** distribution $P(X \in \mathcal{X}, Y \in \mathcal{Y})$
- A loss function: $\ell(y' \in \mathcal{Y}, y \in \mathcal{Y}) \in \mathbb{R}$
- Goal: find a function $h: \mathcal{X} \mapsto \mathcal{Y}$ to minimize
 - $E_{(x,y) \sim P}[\ell(h(x), y)] = E_{(x,y) \sim P}[\ell(\hat{y}, y)]$, where $\hat{y} = h(x)$ is the predicted label

Data pair with true/target label:



More formal definition of supervised learning

- An “unknown” **true** distribution $P(X \in \mathcal{X}, Y \in \mathcal{Y})$
- A loss function: $\ell(y' \in \mathcal{Y}, y \in \mathcal{Y}) \in \mathbb{R}$
- Goal: find a function $h: \mathcal{X} \mapsto \mathcal{Y}$ to minimize
 - $E_{(x,y) \sim P}[\ell(h(x), y)] = E_{(x,y) \sim P}[\ell(\hat{y}, y)]$, where $\hat{y} = h(x)$ is the predicted label
- Training data:
 - $D_{tr} = \{(x_1, y_1), \dots, (x_N, y_N)\}$, where $(x_n \in \mathcal{X}, y_n \in \mathcal{Y}) \sim P$ and $n \in \{1, \dots, N\}$
 - Each pair (x_n, y_n) is IID sampled from P ; $D_{tr} \sim P^N$
 - ML: Apply the learning algorithm to find $h: \mathcal{X} \mapsto \mathcal{Y}$

More formal definition of supervised learning

- An “unknown” **true** distribution $P(X \in \mathcal{X}, Y \in \mathcal{Y})$
- A loss function: $\ell(y' \in \mathcal{Y}, y \in \mathcal{Y}) \in \mathbb{R}$
- Goal: find a function $h: \mathcal{X} \mapsto \mathcal{Y}$ to minimize
 - $E_{(x,y) \sim P}[\ell(h(x), y)] = E_{(x,y) \sim P}[\ell(\hat{y}, y)]$, where $\hat{y} = h(x)$ is the predicted label
- Test data:
 - $D_{te} = \{(x_1, y_1), \dots, (x_M, y_M)\}$, where $(x_m \in \mathcal{X}, y_m \in \mathcal{Y}) \sim P$
 - Each pair (x_m, y_m) is IID sampled from P ; $D_{te} \sim P^M$
 - Compute $\sum_{m=1}^M \ell(h(x_m), y_m)$

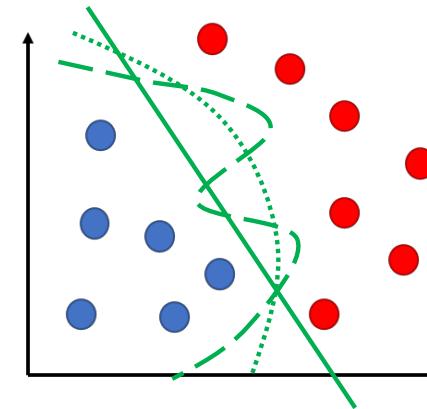
Special cases

- Data distributions:
 - $P(X, Y) = P(X|Y)P(Y) = P(X)P(Y|X)$
 - In some cases, we may assume that there exists an oracle (or labeling or target) function $f: \mathcal{X} \mapsto \mathcal{Y}$ that deterministically generates the true label.
 - $(x, y) \sim P(X)P(Y|X)$ becomes $x \sim P(X)$ and then $y = f(x)$
- Notations:
 - Sometimes people use \mathcal{D} instead of P to denote a distribution

Hypothesis class (model type)

- Instead of allowing arbitrary h , we put some restrictions such that $h \in \mathcal{H}$
- \mathcal{H} : hypothesis class (a set of hypotheses h)

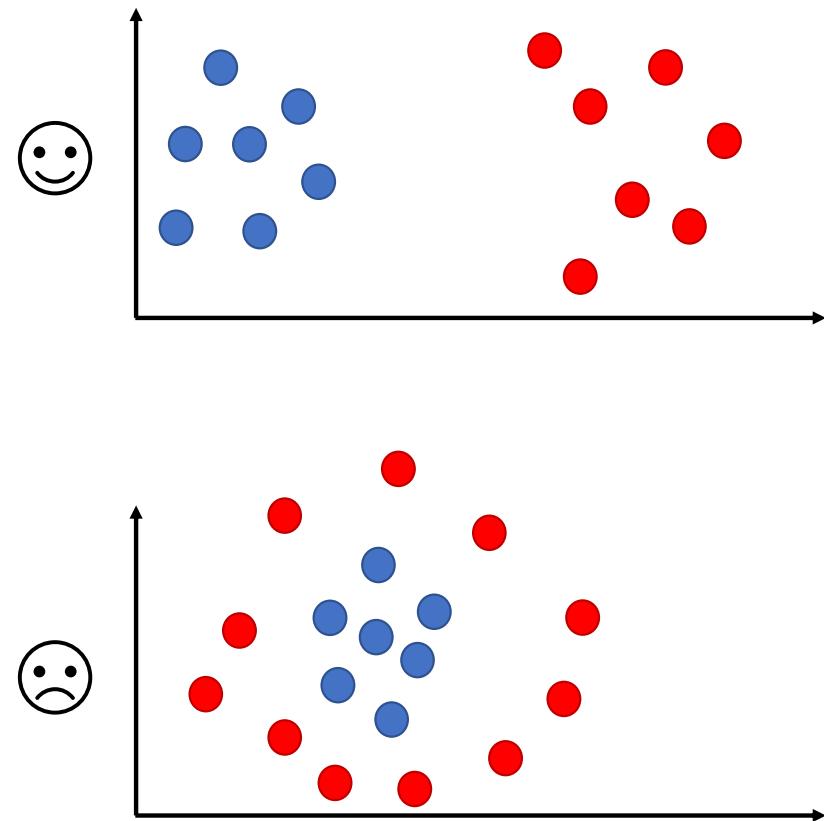
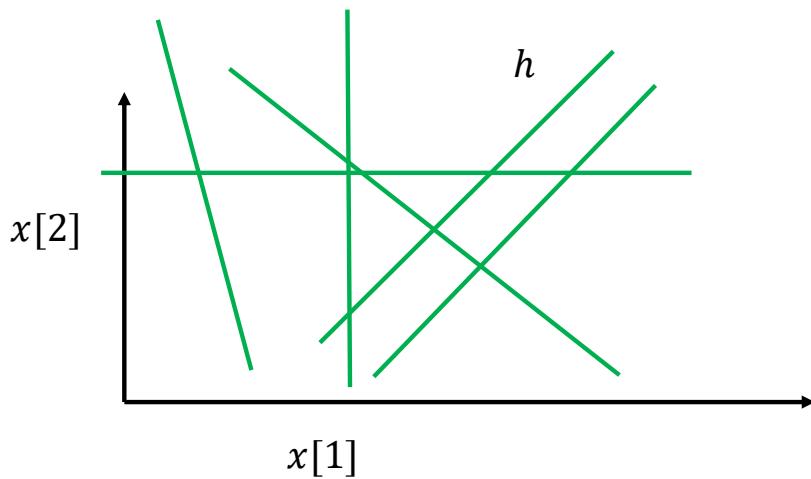
- Linear classifiers/boundaries
- Nonlinear classifiers/boundaries
- Neural networks



- A hypothesis class encodes important *assumptions/prior knowledge/inductive bias* about the *problem (the distribution P)* we are trying to learn
- The **No free lunch theorem** states that every successful ML algorithm must make assumptions
 - There is no single ML algorithm that works for every problem (data distributions)

Hypothesis class (model type)

- Linear classifiers \mathcal{H}



Loss functions ℓ

- For classification
 - 0-1 loss: $\ell(y, y') = 1[y \neq y']$; $1[\text{True}] = 1, 1[\text{False}] = 0$ is called an indicator function
- For regression
 - Square loss: $\ell(y, y') = (y - y')^2$
 - Absolute loss: $\ell(y, y') = |y - y'|$
 - Grow linearly
 - More tolerance to noisy data
- There are many other types of loss functions, and you will see them soon

More formal definition of supervised learning

- Generalization loss/error/risk: $\mathcal{L}(h; P) = \epsilon(h; P) = E_{(x,y) \sim P}[\ell(\mathbf{h}(x), y)]$
 - Aka, true error
 - Realizability (not always): there exists a $h \in \mathcal{H}$ such that $\epsilon(h; P) = 0$ for the 0-1 loss
- Test loss: $\mathcal{L}(h; D_{te}) = \hat{\epsilon}(h; D_{te}) = \frac{1}{M} \sum_{m=1}^M \ell(\mathbf{h}(\mathbf{x}_m), y_m)$
- Training loss: $\mathcal{L}(h; D_{tr}) = \hat{\epsilon}(h; D_{tr}) = \frac{1}{N} \sum_{n=1}^N \ell(\mathbf{h}(\mathbf{x}_n), y_n)$
 - Aka, empirical loss, empirical risk
- Since D_{tr} is random, \hat{h} (the one found by the ML algorithm) is random
 - So are $\hat{\epsilon}(\hat{h}; D_{tr})$, $\hat{\epsilon}(\hat{h}; D_{te})$, and $\epsilon(\hat{h}; P)$

More in learning theories!

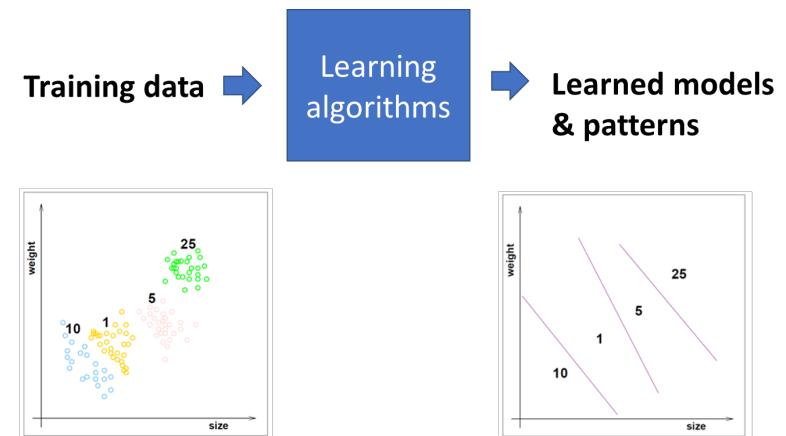
Questions?



THE OHIO STATE UNIVERSITY

ML problem solving given D_{tr} (inner loop)

- Step 1: Select an ML algorithm
 - We think appropriate for the learning problem
 - Include selecting \mathcal{H}
 - More to be defined
- Step 2: Find the “best” $\hat{h} \in \mathcal{H}$ from D_{tr}
 - The “best” depends on the ML algorithm
 - Ideally, “best” for minimizing $\epsilon(\hat{h}; P)$
 - Usually involve an optimization problem



ML problem solving given D_{tr} (outer loop)

- How to pick an ML learning algorithm if there are several choices?
 - Goal: minimize the generalization error or test error
 - Problem: only have the training data
 - Solution: separate D_{tr} into D_{base} and D_{val} ; treat D_{val} as “pseudo” test data
 - Separate “randomly” to keep D_{base}, D_{val} distributionally similar, not simply by first-vs-second half!
- For each learning algorithm (and its \mathcal{H})
 - Find the “best” $\hat{h} \in \mathcal{H}$ from D_{base}
 - Compute $\hat{\epsilon}(\hat{h}; D_{val})$
- Choose the best algorithm whose $\hat{\epsilon}(\hat{h}; D_{val})$ is the smallest
 - Find the “best” $\hat{h} \in \mathcal{H}$ from D_{tr} , using the best algorithm

Empirical risk minimization (ERM)

- A widely used ML algorithm (together with different \mathcal{H})
- Given D_{tr} (or D_{base}), output $\hat{h} \in \mathcal{H}$
 - $\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} \hat{\epsilon}(h; D_{tr}) = \operatorname{argmin}_{h \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^N \ell(\textcolor{red}{h}(x_n), y_n)$
 - “Hope” that it will work well on $\epsilon(\hat{h}; P)$
- **When will the “hope” have some theoretical guarantee?**
 - $D_{tr} \sim P^N$ (or from a similar distribution: e.g., train on synthetic data, test on real data)
 - \mathcal{H} is with appropriate assumptions and an appropriate set size
 - N is large enough

Generalization

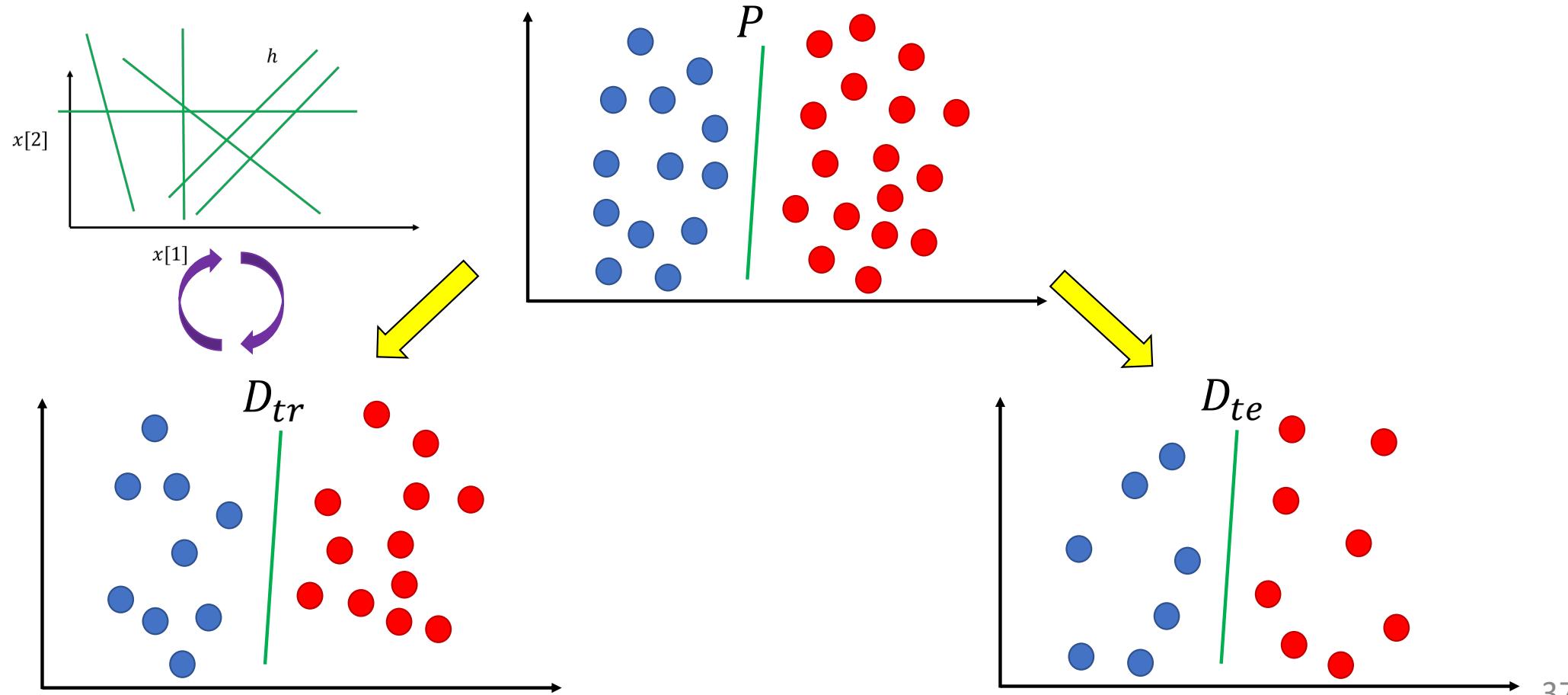
- Can we simply choose the ML algorithm whose $\hat{\epsilon}(\hat{h}; D_{tr})$ is minimized?
 - Example: $\hat{h}(\mathbf{x}) = \begin{cases} y_n, & \text{if } \exists (\mathbf{x}_n, y_n) \in D_{tr} \text{ s.t. } \mathbf{x} = \mathbf{x}_n \\ 0, & \text{otherwise} \end{cases}$
 - $\hat{\epsilon}(\hat{h}; D_{tr}) = 0$, but know nothing about data beyond D_{tr}
- How to resolve this?
 - We can detect this problem by finding $\hat{h}(\mathbf{x})$ from D_{base} and evaluating $\hat{\epsilon}(\hat{h}; D_{val})$
 - This is why we separate D_{tr} into D_{base} and D_{val} to choose an ML algorithm

Questions?



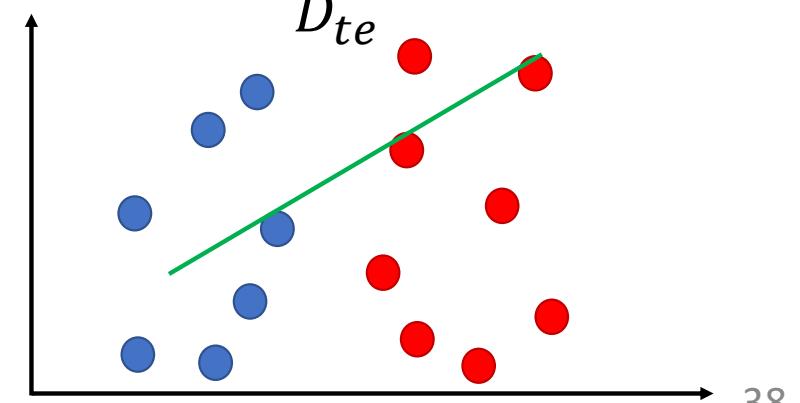
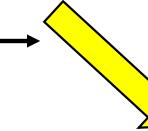
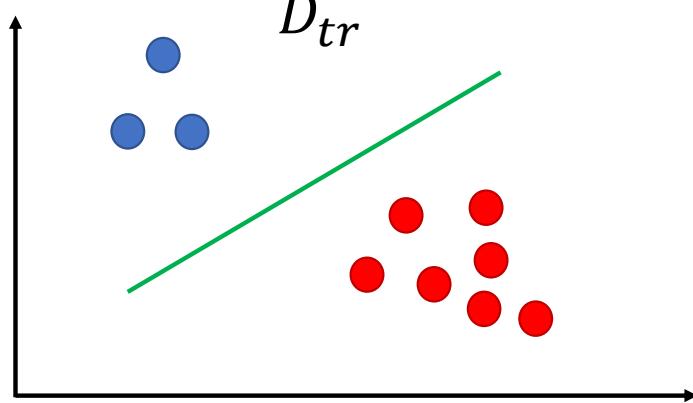
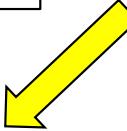
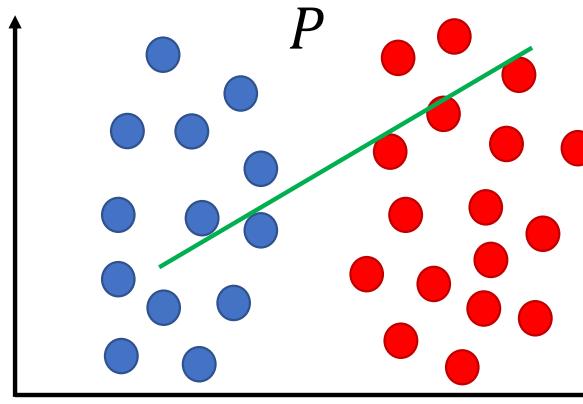
THE OHIO STATE UNIVERSITY

Why does data size matter?



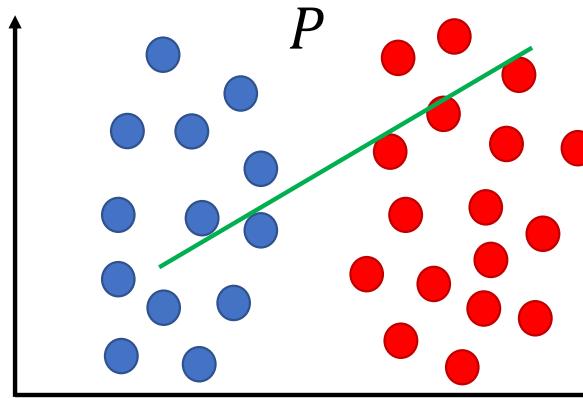
Why does data size matter?

D_{tr} not
representative of P
to find a good \hat{h}

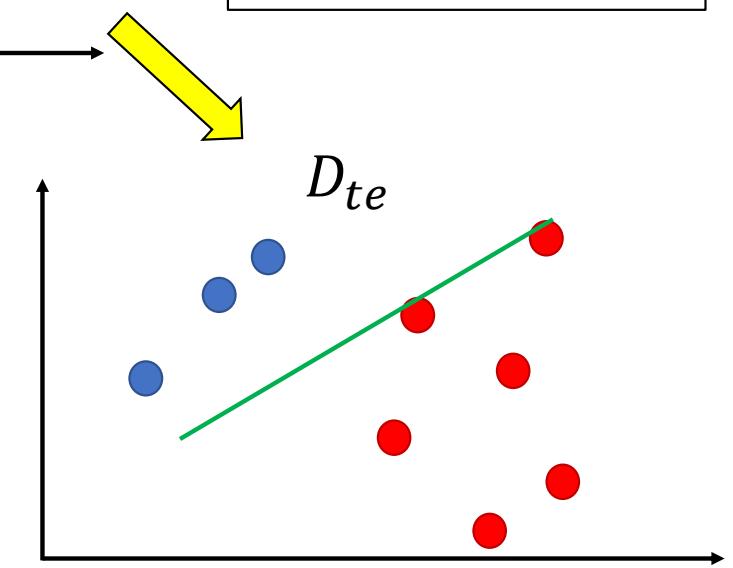
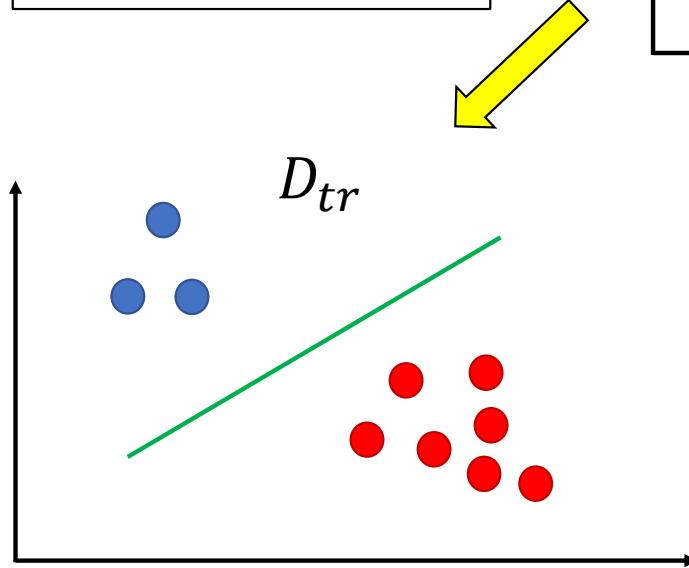


Why does data size matter?

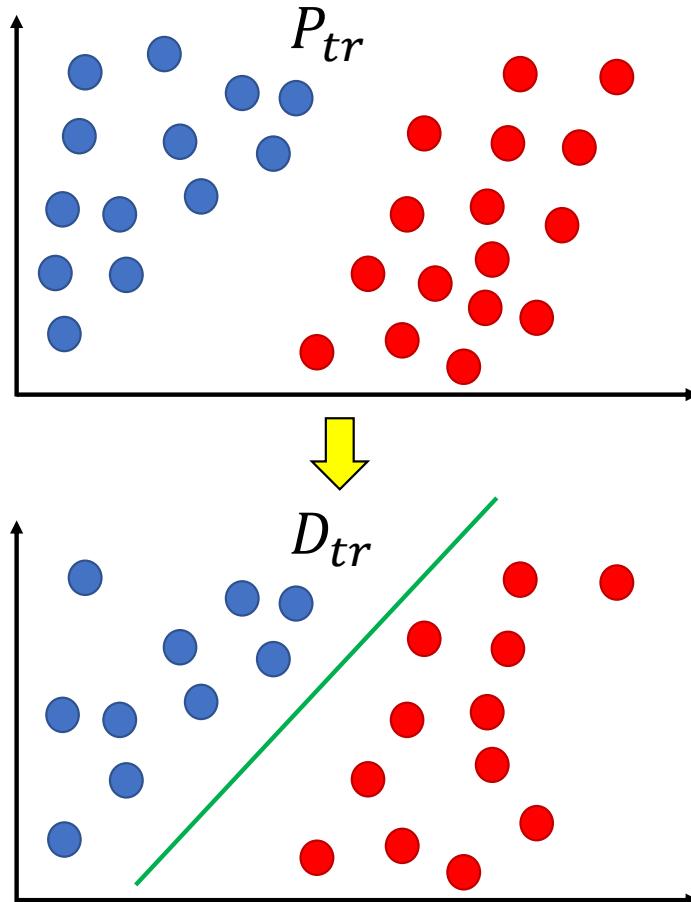
D_{tr} not
representative of P
to find a good \hat{h}



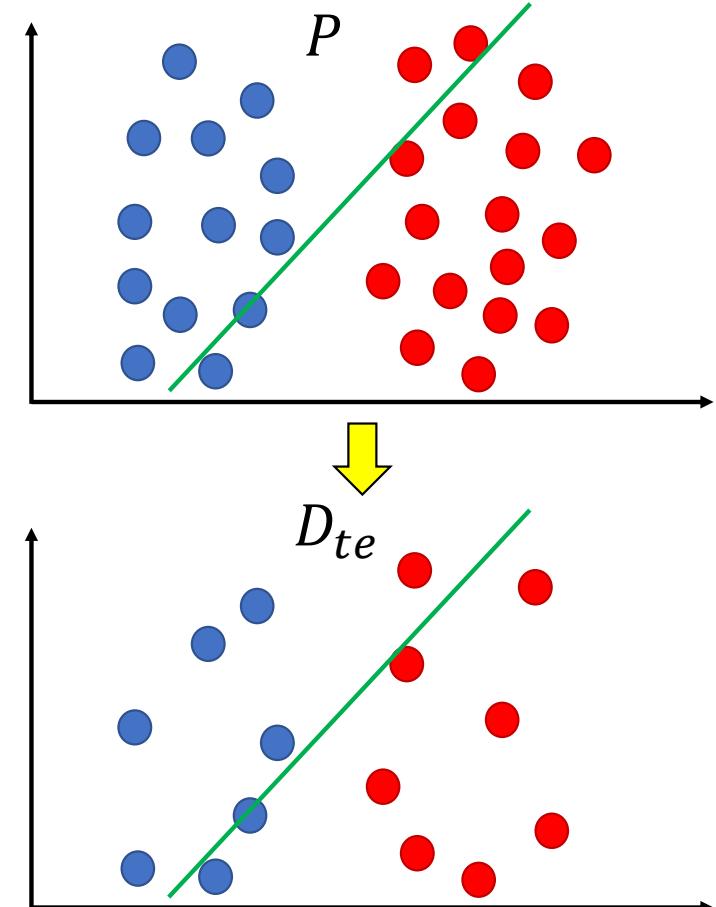
D_{te} not
representative of P
to evaluate \hat{h}



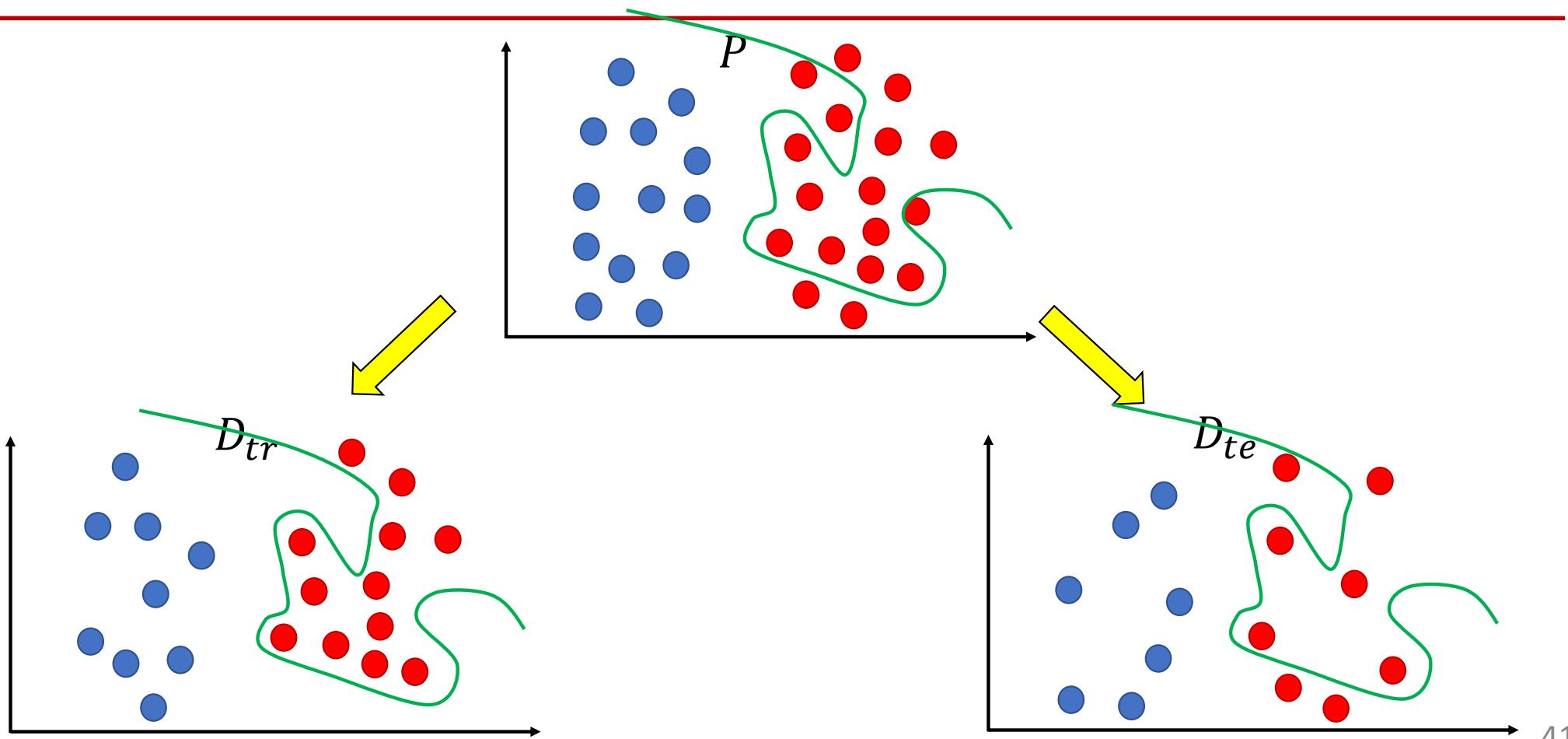
Why do the similar data distributions matter?



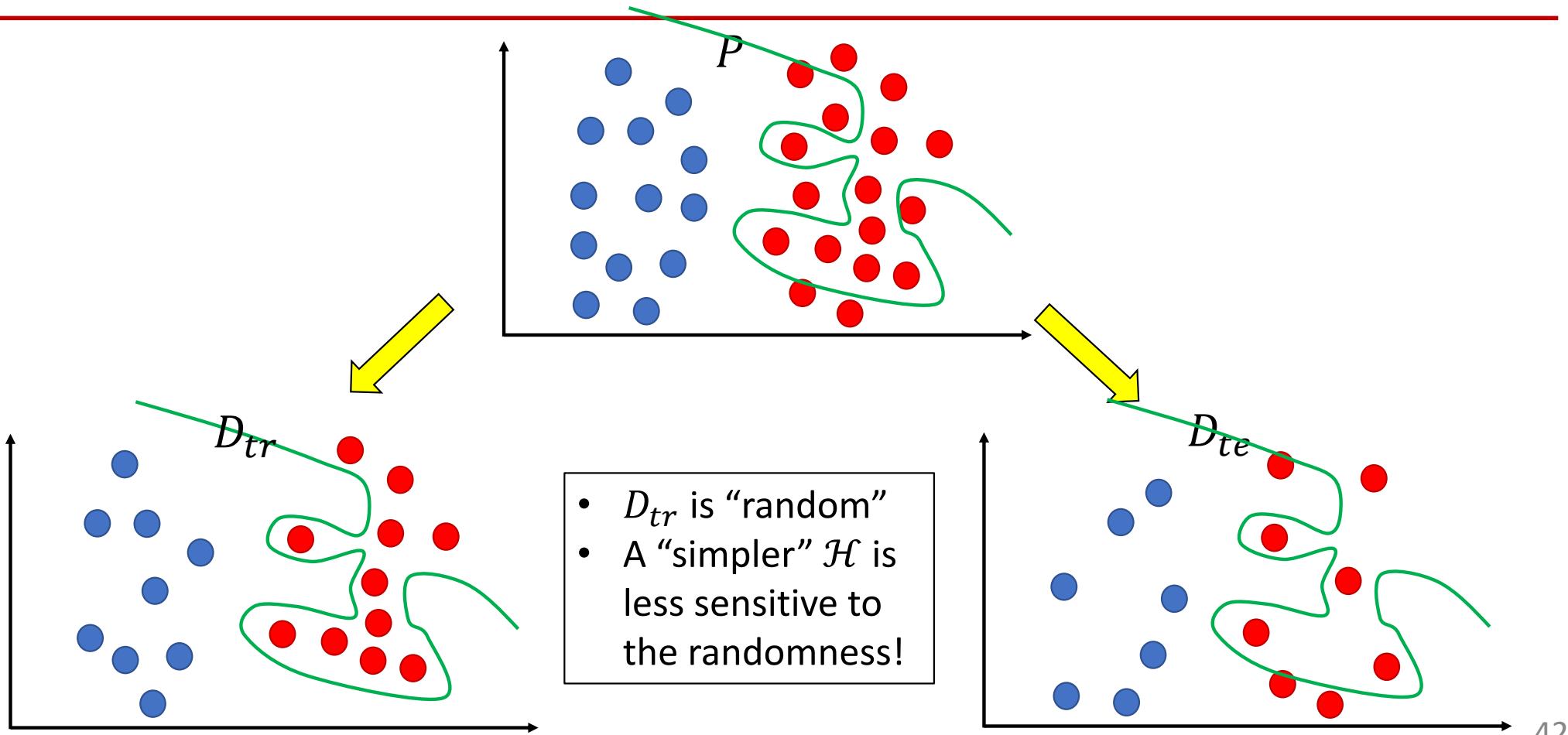
D_{tr} not
representative of P
to find a good \hat{h}



Why does the hypothesis class matter?



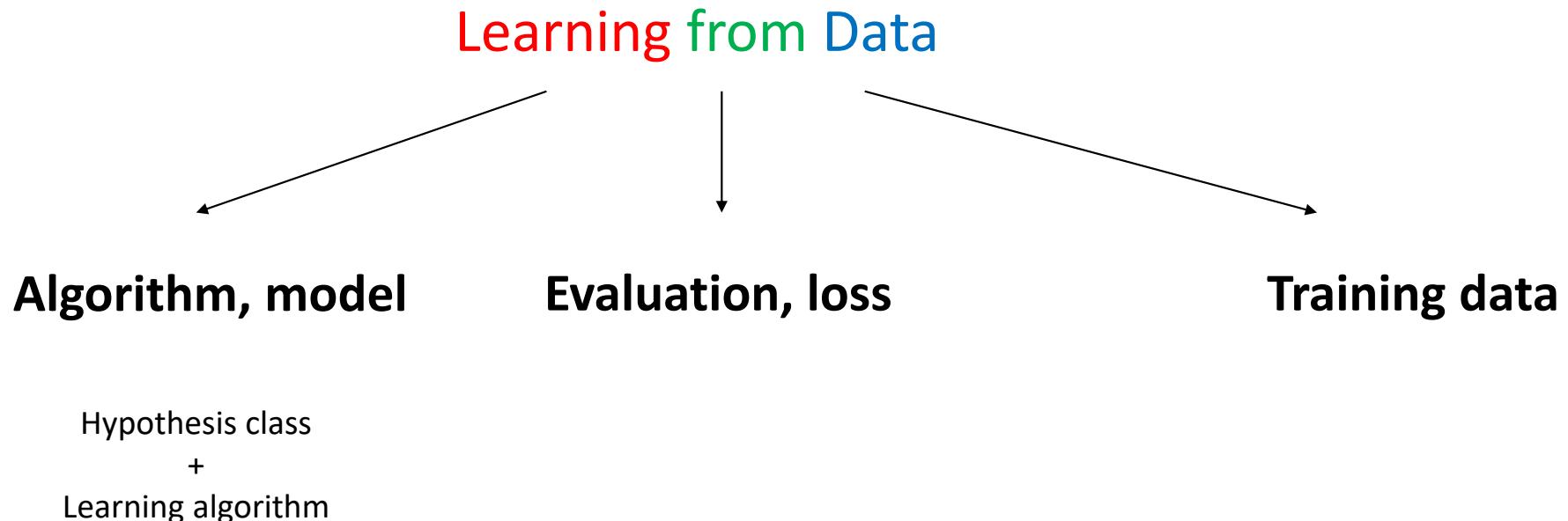
Why does the hypothesis class matter?



More

- Learning theories
- Concentration theory, Hoeffding's inequality, VC dimension, ...

Summary



Why do we need linear algebra?

- Machine learning:

- Learning patterns from data
- A data instance can usually be represented by a vector
- A dataset can usually be represented by a matrix
- The “pattern” can usually be represented by a vector or a matrix!
 - A coin with weight * 0.2 + diameter * 0.04 < 2 \Rightarrow 10 cent
- To learn the right patterns, we need computation between possible patterns and data!



weight
diameter $\begin{bmatrix} 5.6 \\ 23.5 \end{bmatrix}$

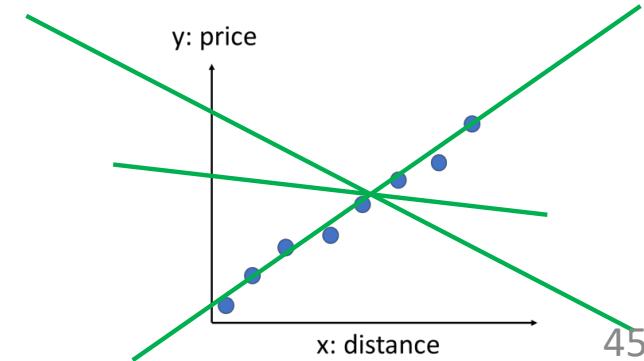


$$2 - \begin{bmatrix} 0.2 \\ 0.04 \end{bmatrix}^T \times \begin{bmatrix} 5.6 \\ 23.5 \end{bmatrix}$$

- Linear algebra:

- Powerful for computation on vectors and matrices
- More important to get the meaning of the computation!
- Please review the linear algebra slides!

Caution! Python indices start with 0!



TODO for math/programming review

- See the beginning slides and the course website for suggested reading

CSE 5523: Nearest Neighbor



THE OHIO STATE UNIVERSITY

Important for this week

- **Register:** If you were on the waitlist but have not been enrolled, please talk to me after the class
- **Register for the class on piazza:** our platform for discussion and communication
 - https://piazza.com/osu/autumn2024/cse5523_chao (please use name.#@osu.edu)
 - Access code: osu-cse-5523-AU24-chao
- **Math review/self-diagnostic:** do Homework #0 and suggested materials on the website --- **extremely important to check your readiness for the course**
- **Decision:** stay or drop
- **Office hours:** start this Thursday
 - Mine: Tuesday 11 am - 12 pm, Thursday 6 pm - 7 pm (DL 587)
 - TA (Tai-Yu Pan): Monday 11 am - 12 pm, Friday 4 pm - 5 pm (BE406, #12 and #13)

Interest in the Buckeye AutoDrive Team?

Interested in the growing field of Autonomous Vehicles (AV)?
Want to transform a Chevy Bolt into a fully autonomous vehicle?



Website

We have opportunities for CSE,
MAE and ECE Students!

Gain hands-on technical skills while
also networking with industry
professionals at SAE & GM!



Interest Survey



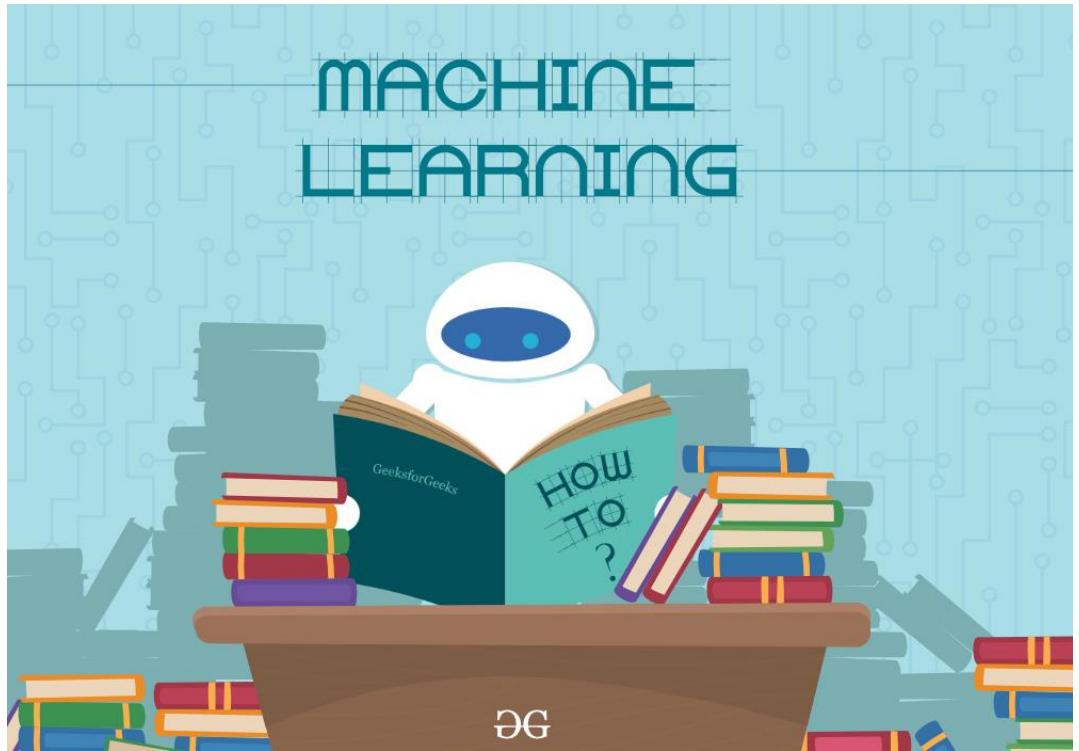
Today

Review

- ML problem solving
- Training, validation, testing

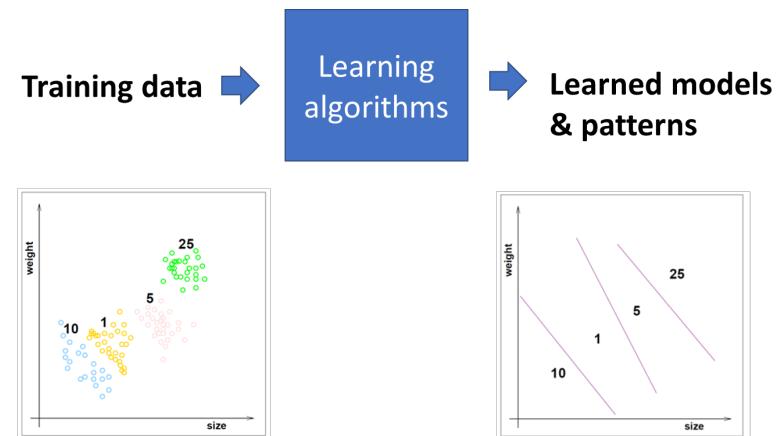
Nearest neighbor

- Basic and variants
- Properties
- Curse of dimensionality
- Bayes optimal classifier



ML problem solving given D_{tr} (inner loop)

- Step 1: Select an ML algorithm
 - We think appropriate for the learning problem
 - Include selecting hypothesis class \mathcal{H}
 - More to be defined
- Step 2: Find the “best” $\hat{h} \in \mathcal{H}$ based on D_{tr}
 - The “best” depends on the ML algorithm
 - Ideally, “best” for minimizing $\epsilon(\hat{h}; P)$
 - Usually involve an optimization problem



ML problem solving given D_{tr} (outer loop)

- How to pick an ML learning algorithm if there are several choices?
 - Goal: minimize the generalization error or test error
 - Problem: only have the training data
 - Solution: separate D_{tr} into D_{base} and D_{val} ; treat D_{val} as “pseudo” test data
 - Separate “randomly” to keep D_{base}, D_{val} distributionally similar, not simply by first-vs-second half!
- For each learning algorithm (and its \mathcal{H})
 - Find the “best” $\hat{h} \in \mathcal{H}$ based on D_{base}
 - Compute $\hat{\epsilon}(\hat{h}; D_{val})$
- Choose the best algorithm whose $\hat{\epsilon}(\hat{h}; D_{val})$ is the smallest
 - Find the “best” $\hat{h} \in \mathcal{H}$ based on D_{tr}

Empirical risk minimization (ERM)

- A widely used ML algorithm (together with different \mathcal{H})
- Given D_{tr} (or D_{base}), output $\hat{h} \in \mathcal{H}$
 - $\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} \hat{\epsilon}(h; D_{tr}) = \operatorname{argmin}_{h \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^N \ell(\textcolor{red}{h}(x_n), y_n)$
 - “Hope” that it will work well on $\epsilon(\hat{h}; P)$
- When will the “hope” have some theoretical guarantee?
 - $D_{tr} \sim P^n$ (or from a similar distribution: e.g., train on synthetic data, test on real data)
 - \mathcal{H} is with appropriate “assumptions” and an appropriate set size
 - N is large enough

Questions?



THE OHIO STATE UNIVERSITY

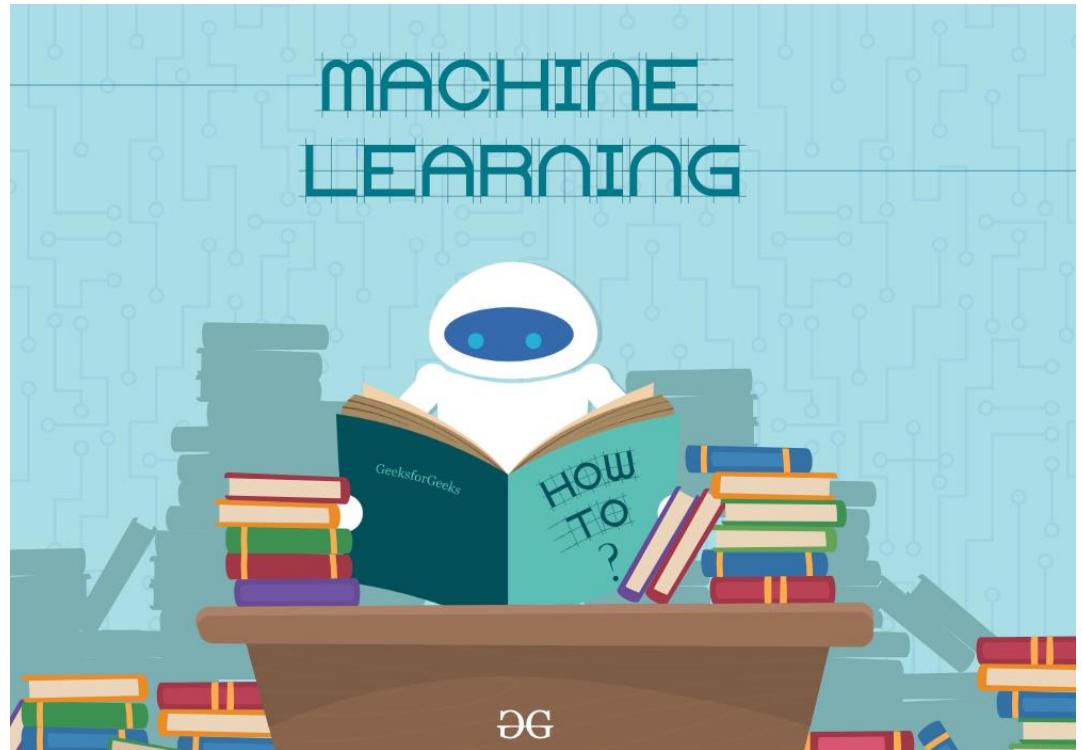
Today

Review

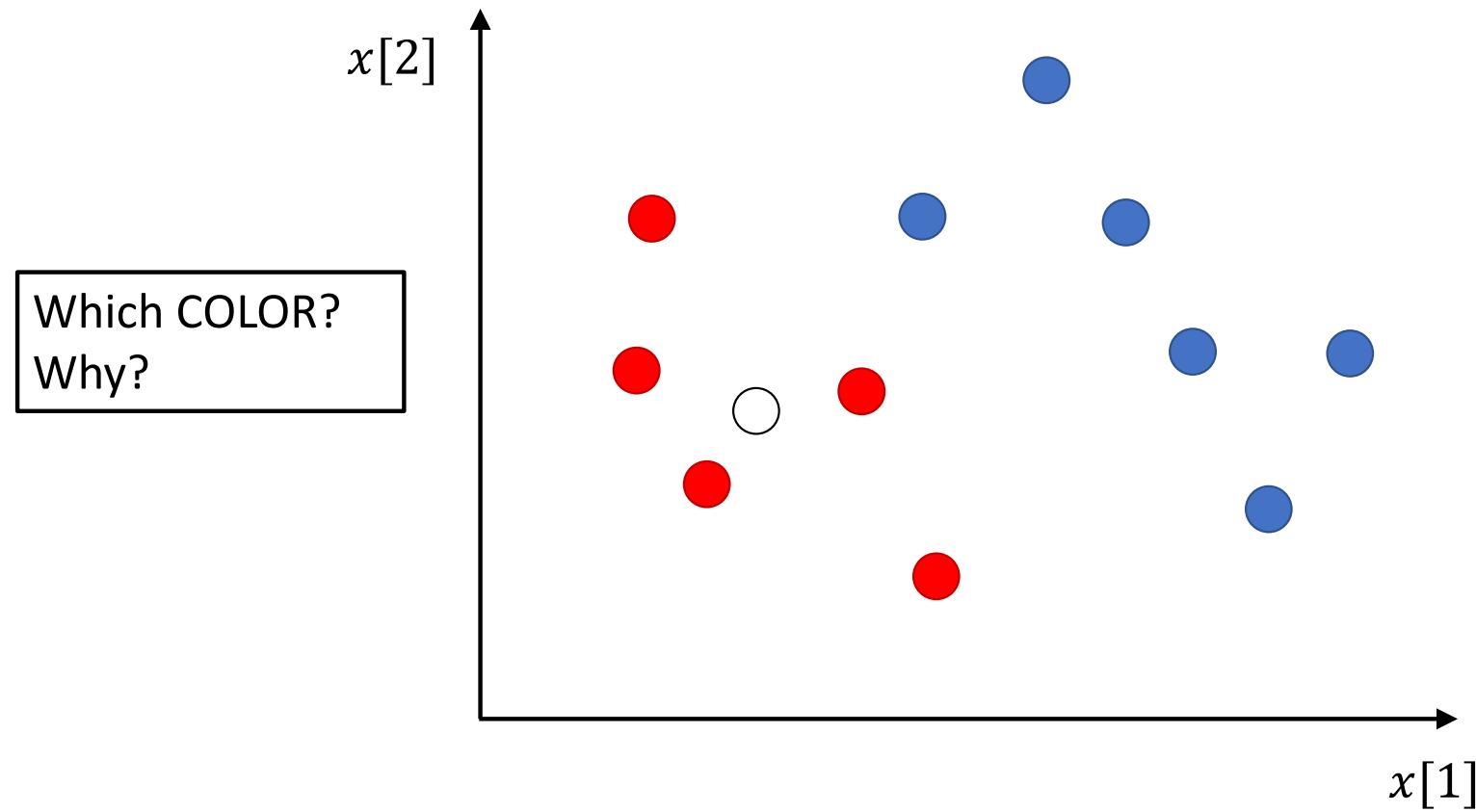
- ML problem solving
- Training, validation, testing

Nearest neighbor

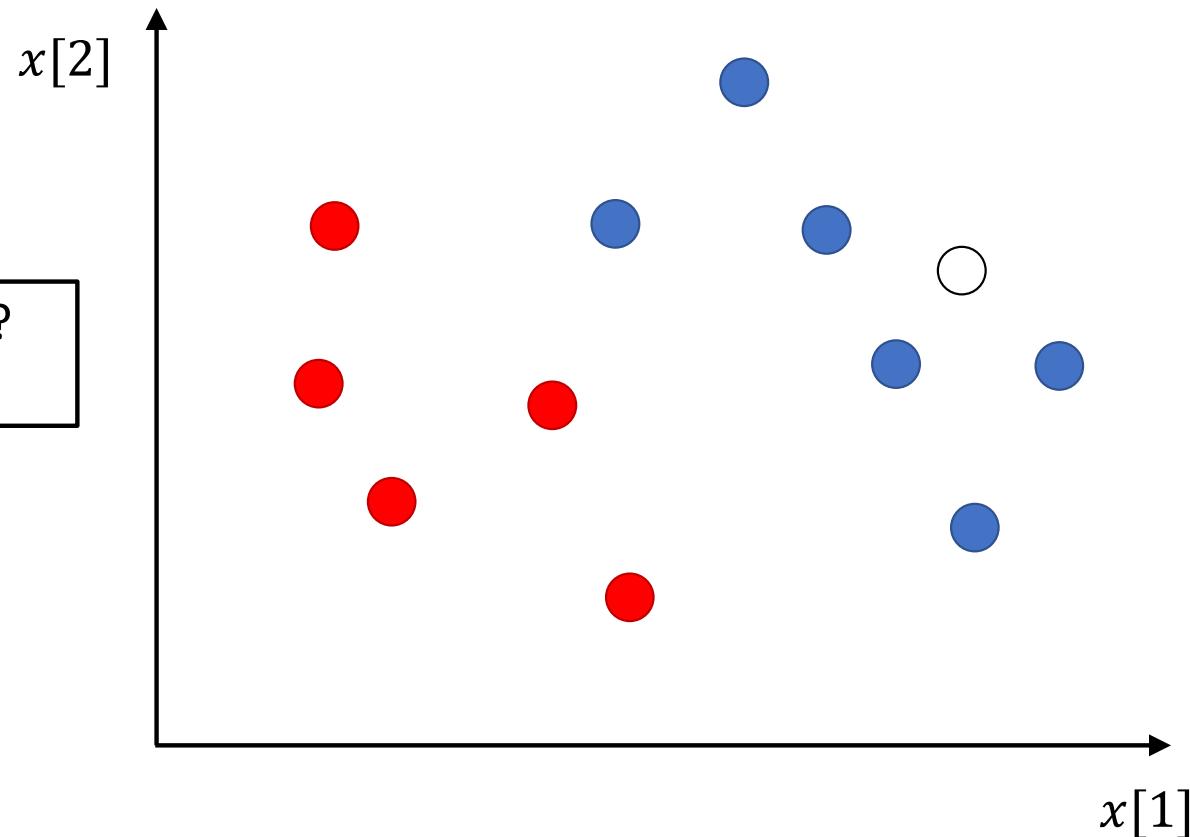
- Basic and variants
- Properties
- Curse of dimensionality
- Bayes optimal classifier



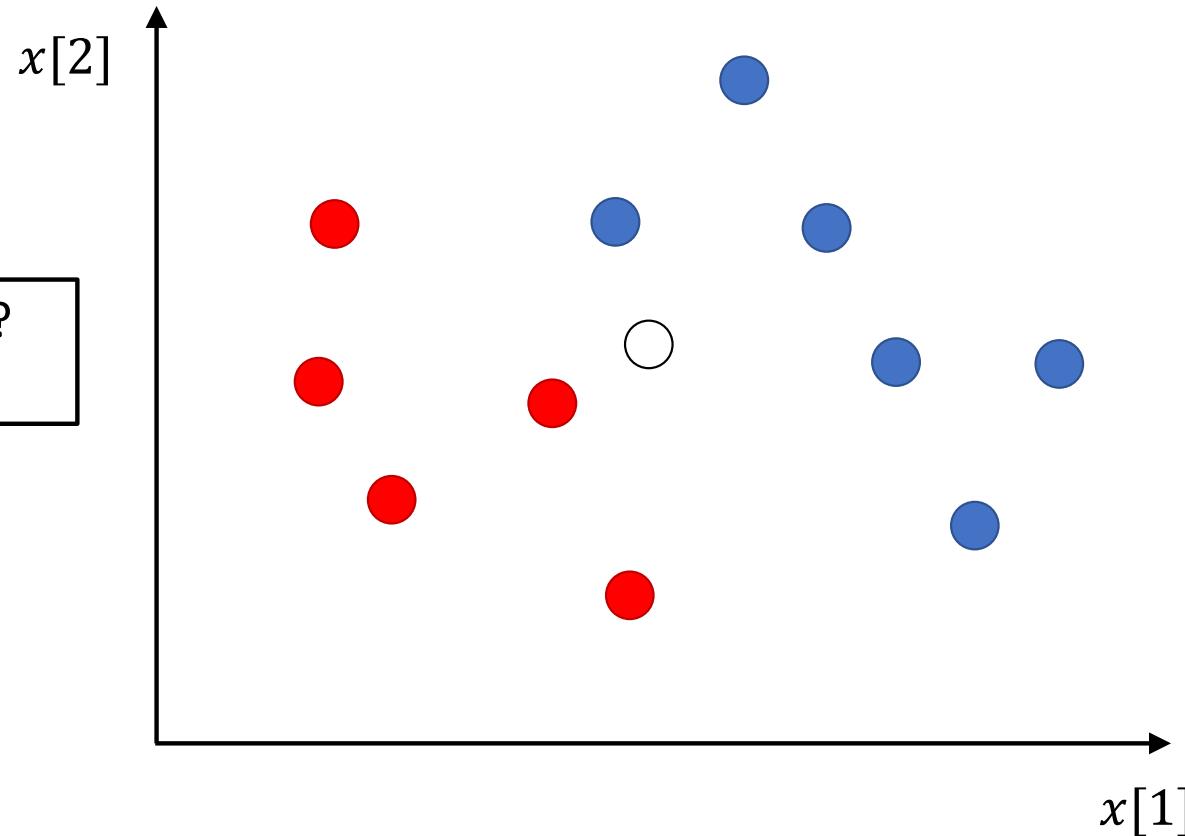
Nearest neighbor



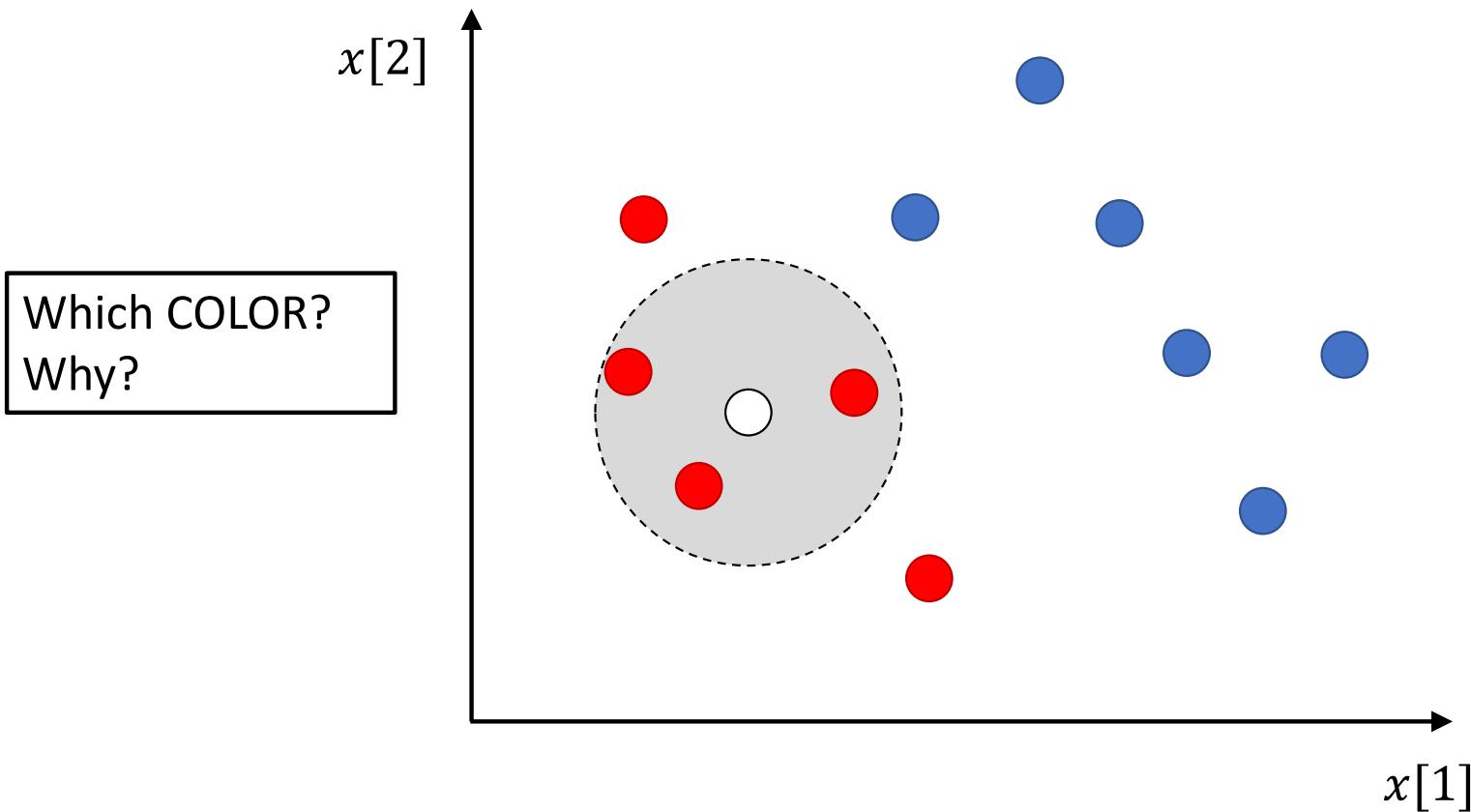
Nearest neighbor



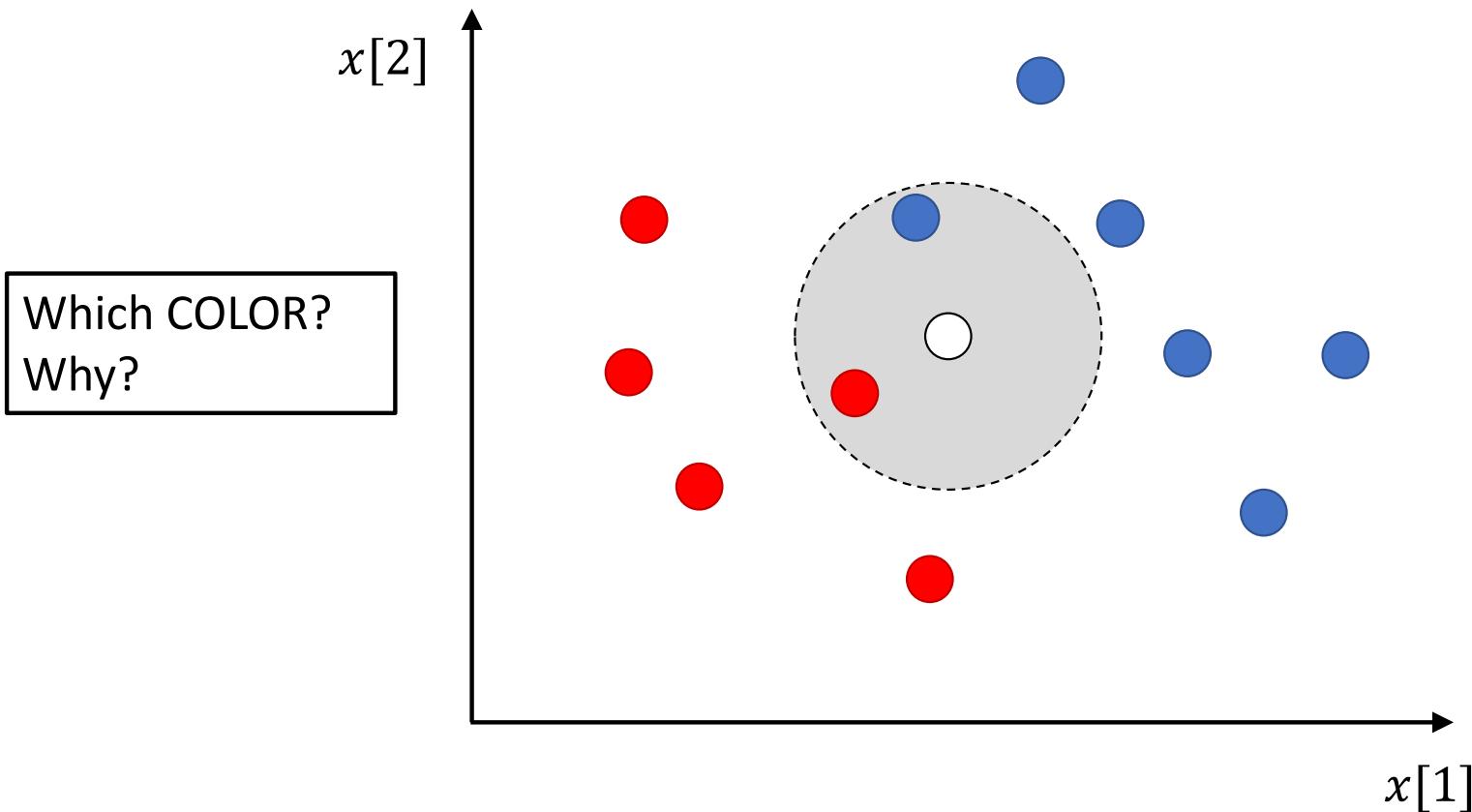
Nearest neighbor



Nearest neighbor



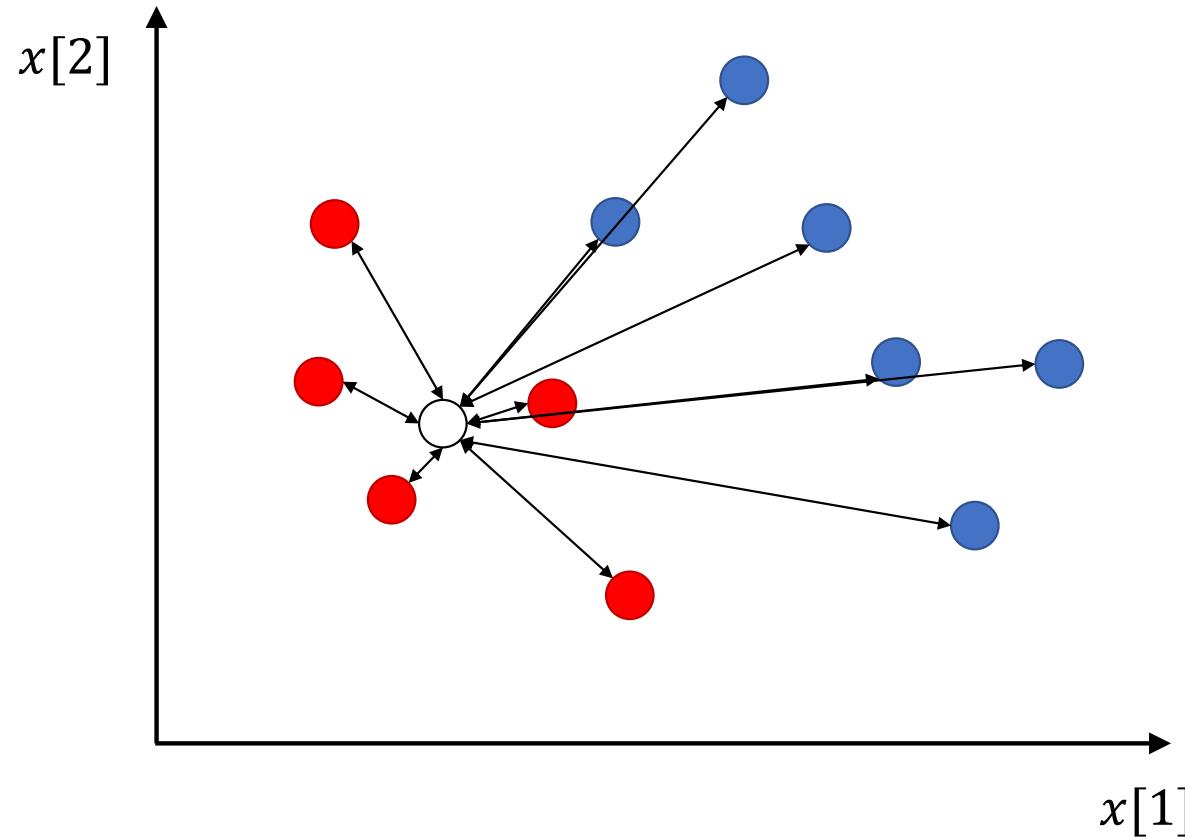
Nearest neighbor



Nearest neighbor

- What is the assumption?
 - Similar Inputs (feature vectors) have similar outputs (labels)
- What does it mean by **similar**?
 - Smaller **distance**, larger **similarity**
 - Larger **distance**, smaller **similarity**
- Classification rule:
 - For a test data input x , assign the label of its closest training data instance in D_{tr}

Nearest neighbor



Nearest neighbor

- **Given:**
 - A distance measure $\text{dist}(\cdot, \cdot)$
 - $D_{tr} = \{(x_1, y_1), \dots, (x_N, y_N)\}$, $t = \infty$ (min distance)

- **Input:** a test data input x

- **NN classification rule (implementation):**

$t = \infty$

$n = 1 : N$

if $\text{dist}(x, x_n) < t$

$\hat{y} = y_n$

$t = \text{dis}(x, x_n)$

Return: \hat{y}

Distance

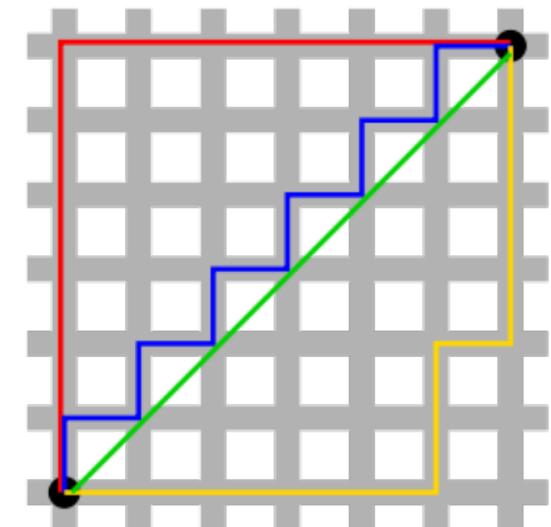
- Nearest neighbor classifiers rely on a distance metric
- The better that metric reflects label similarity, the better the classifier is.

○ L_2 distance (Euclidean): $\|x - x_n\|_2 = \left(\sum_{d=1}^D |x[d] - x_n[d]|^2 \right)^{1/2}$

○ L_1 distance: $\|x - x_n\|_1 = \sum_{d=1}^D |x[d] - x_n[d]|$

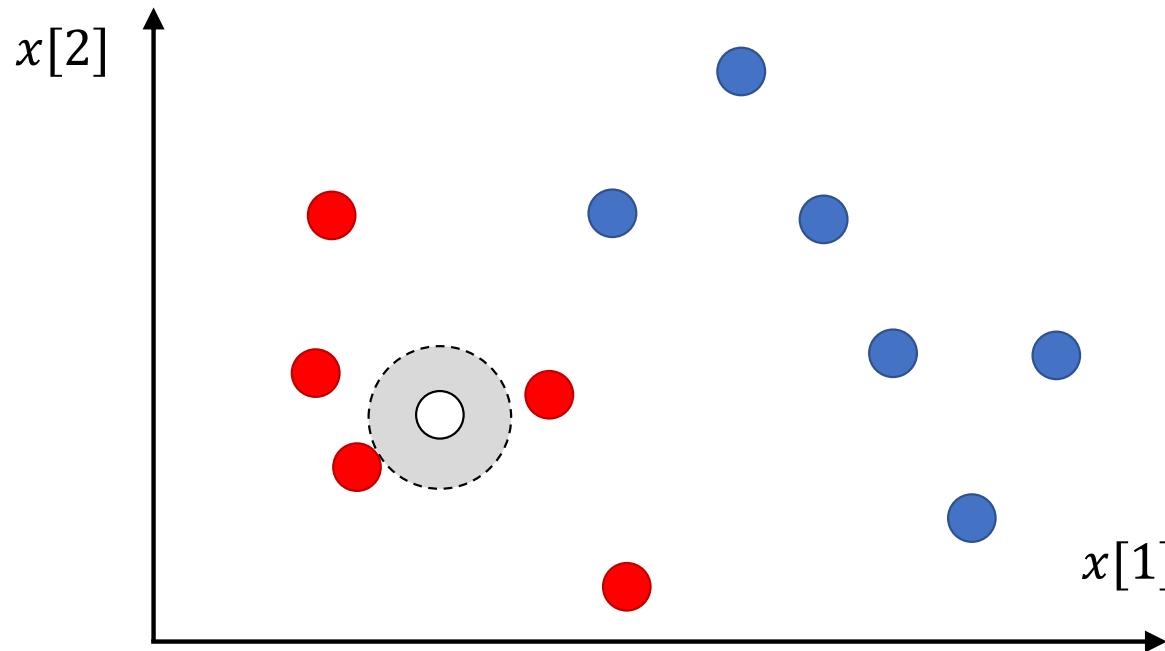
○ L_p norm: $\|x - x_n\|_p = \left(\sum_{d=1}^D |x[d] - x_n[d]|^p \right)^{1/p}$

x : 2-dimensional vectors



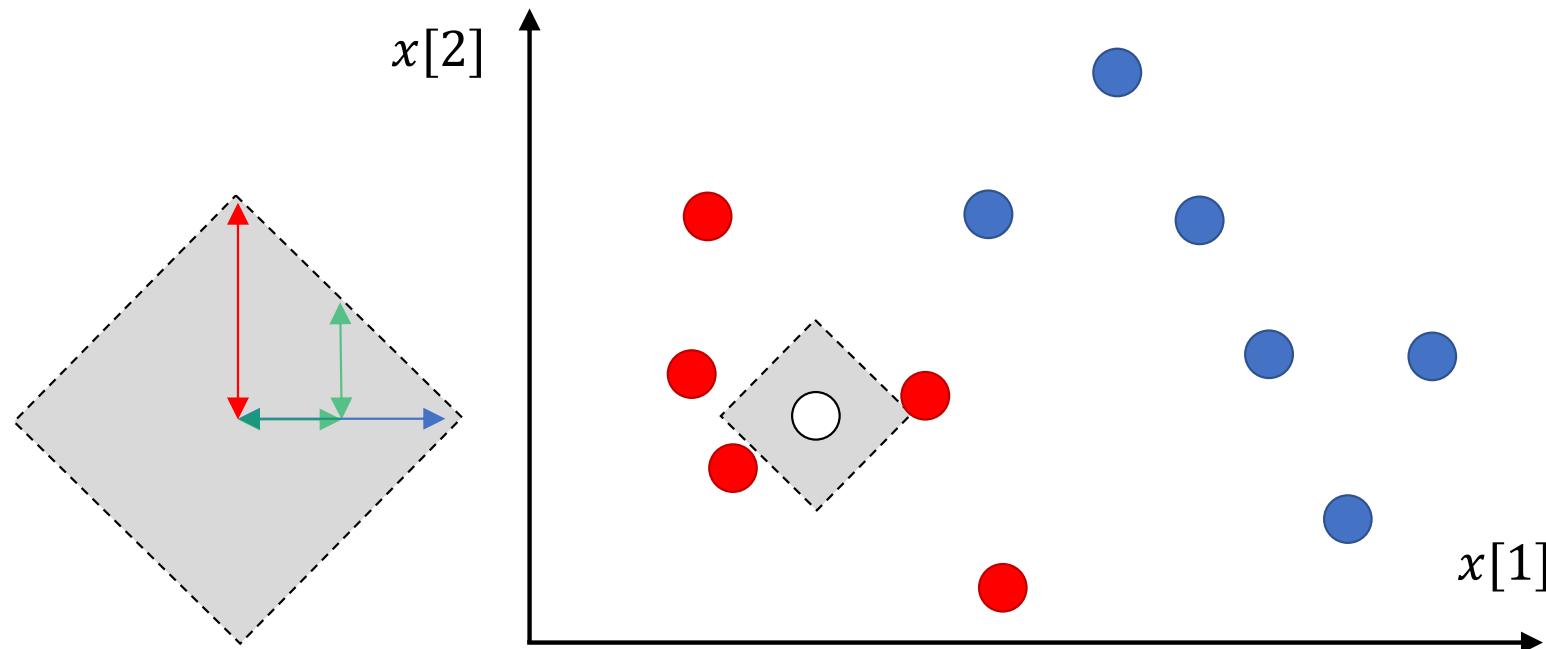
Green line is Euclidean distance.
Red, Blue, and Yellow lines are L_1 distance

Nearest neighbor (Euclidean distance)



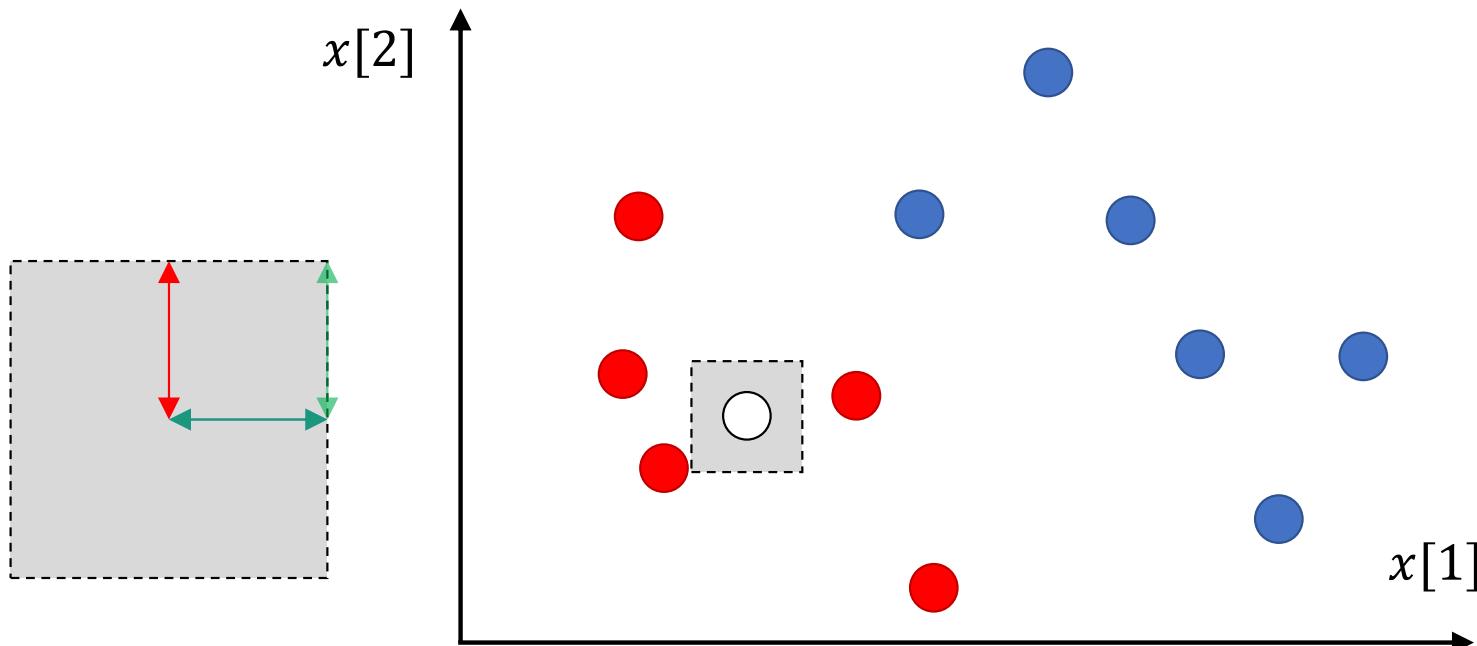
$$\|x - x_n\|_2 = \left(\sum_{d=1}^D |x[d] - x_n[d]|^2 \right)^{1/2}$$

Nearest neighbor (Manhattan distance)



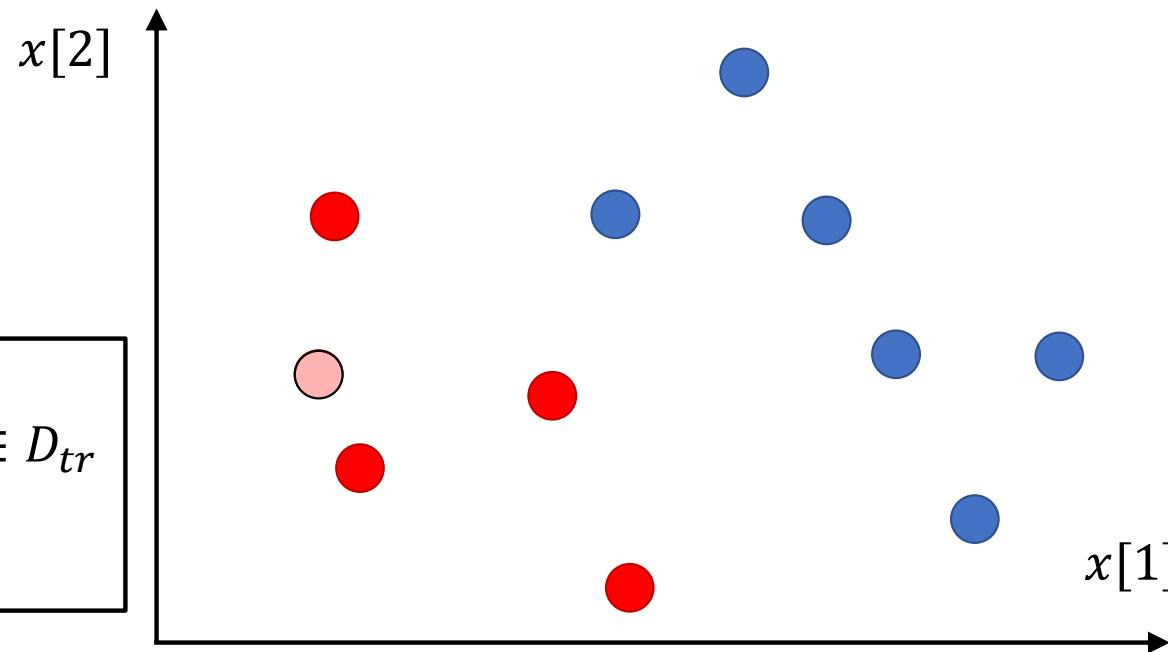
$$\|x - x_n\|_1 = \left(\sum_{d=1}^D |x[d] - x_n[d]|^1 \right)^{1/1}$$

Nearest neighbor



$$\|x - x_n\|_{\infty} = \left(\sum_{d=1}^D |x[d] - x_n[d]|^{\infty} \right)^{1/\infty} \approx \max_d |x[d] - x_n[d]|$$

Nearest neighbor (Euclidean distance)



$$\|x - x_n\|_2 = \left(\sum_{d=1}^D |x[d] - x_n[d]|^2 \right)^{1/2}$$

Questions?

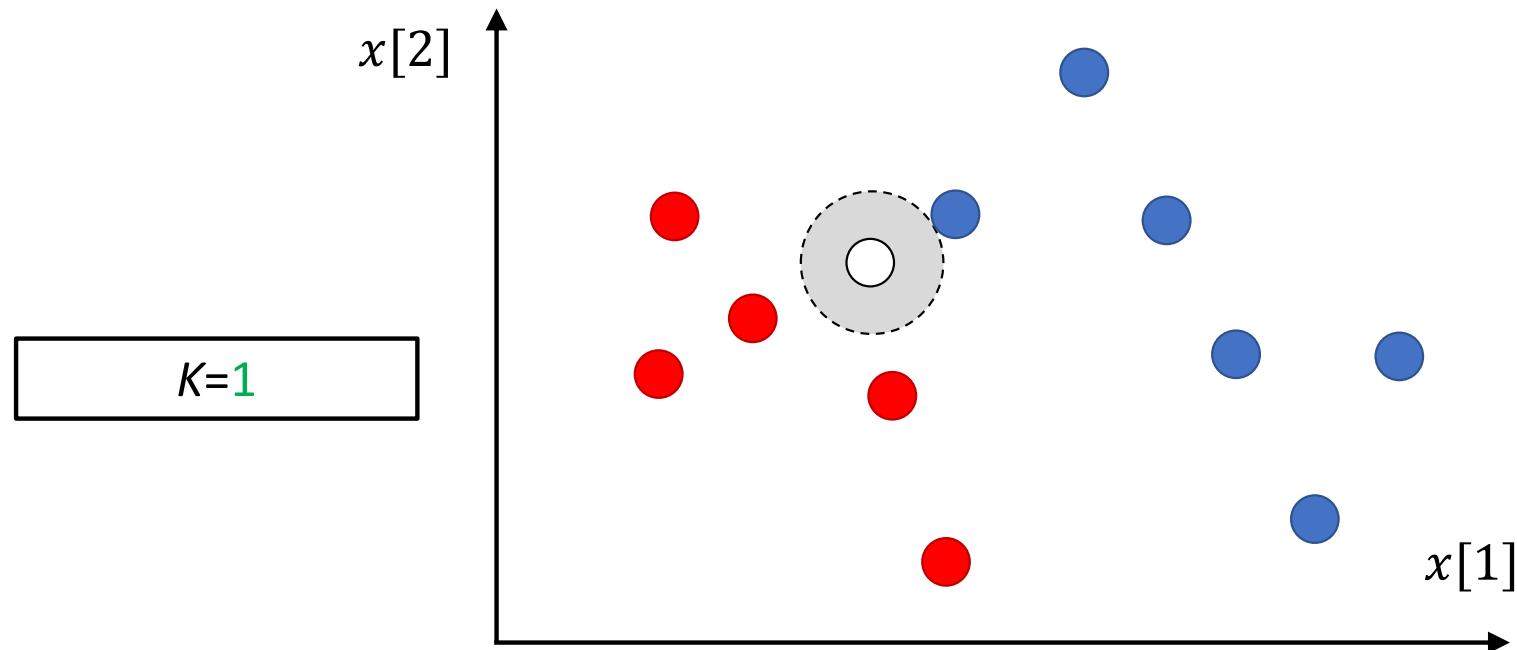


THE OHIO STATE UNIVERSITY

K-Nearest neighbors (KNN)

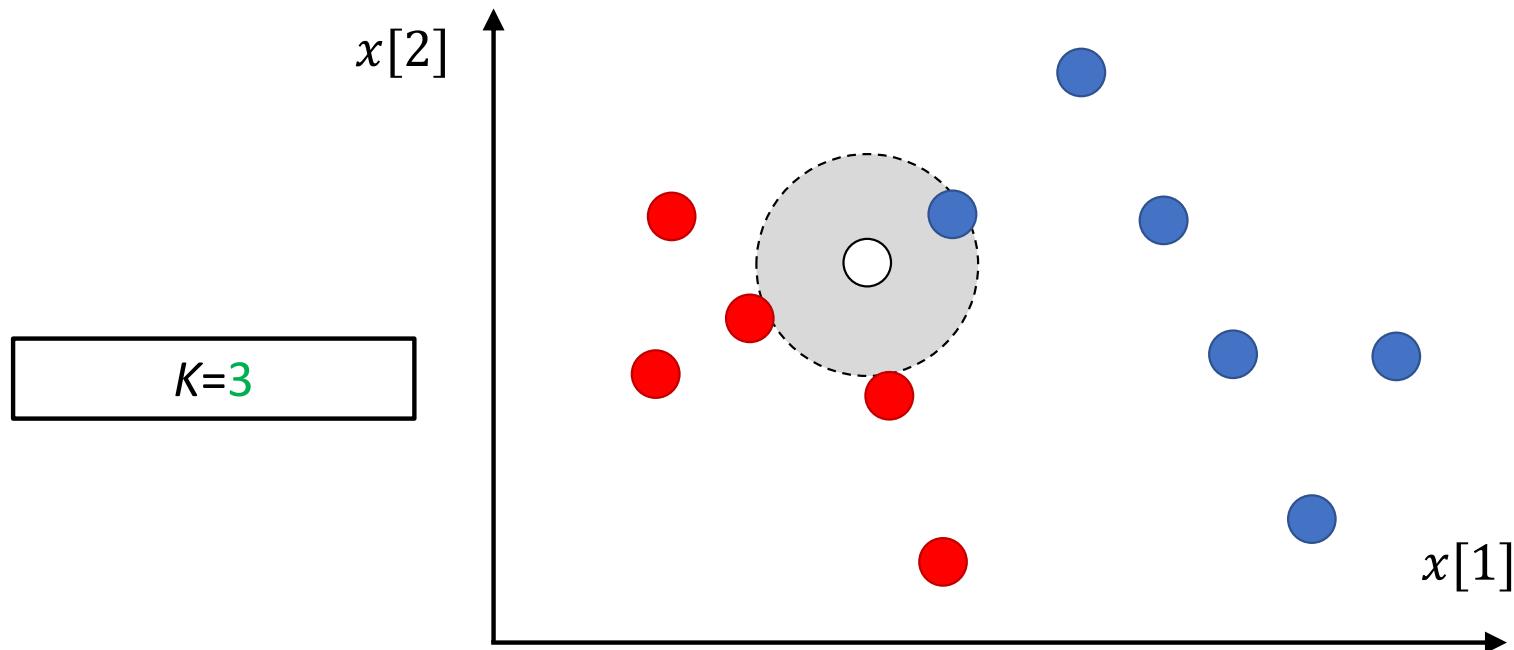
- What is the assumption?
 - Similar Inputs (feature vectors) have similar outputs (labels)
- Classification rule:
 - For a test data input x , assign the **most common label** amongst its K **most similar** training data instances in D_{tr}
 - aka, majority vote

KNN (Euclidean distance)



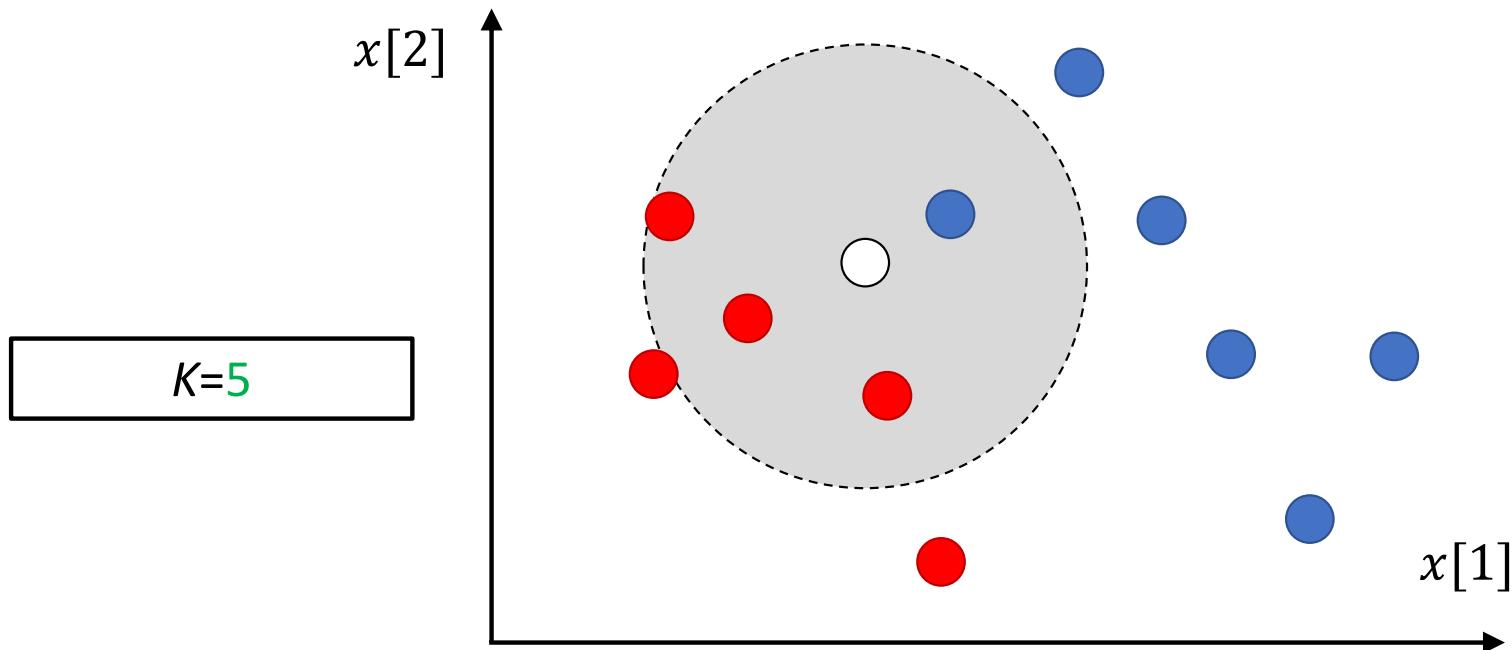
$$\|x - x_n\|_2 = \left(\sum_{d=1}^D |x[d] - x_n[d]|^2 \right)^{1/2}$$

KNN (Euclidean distance)



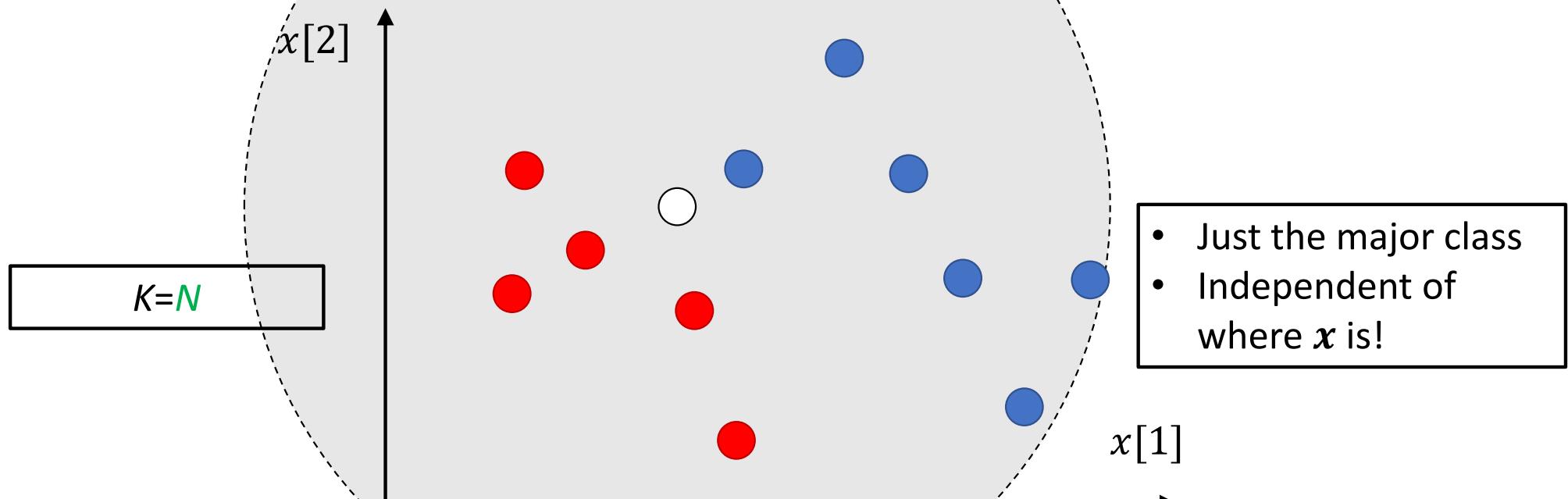
$$\|x - x_n\|_2 = \left(\sum_{d=1}^D |x[d] - x_n[d]|^2 \right)^{1/2}$$

KNN (Euclidean distance)



$$\|x - x_n\|_2 = \left(\sum_{d=1}^D |x[d] - x_n[d]|^2 \right)^{1/2}$$

KNN (Euclidean distance)



$$\|x - x_n\|_2 = \left(\sum_{d=1}^D |x[d] - x_n[d]|^2 \right)^{1/2}$$

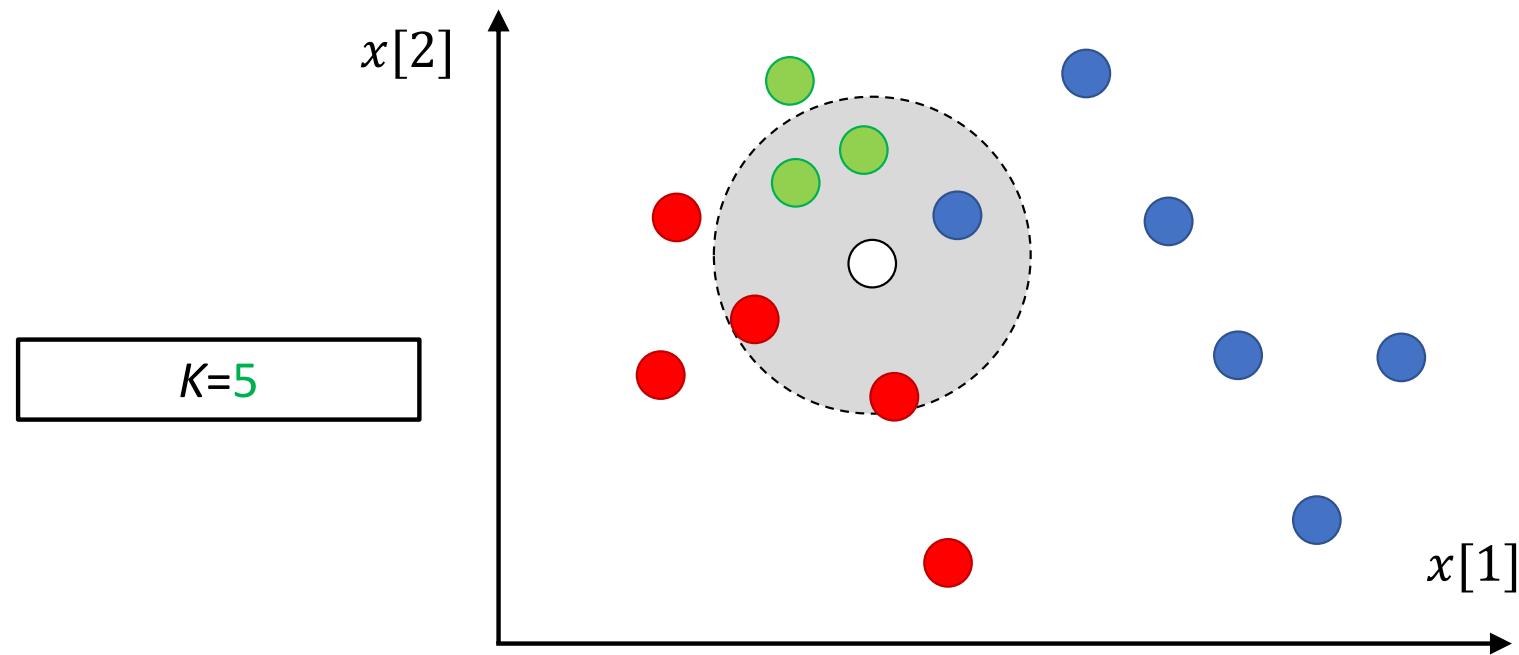
KNN

- **Given:**
 - A distance measure $\text{dist}(\cdot, \cdot)$
 - $D_{tr} = \{(x_1, y_1), \dots, (x_N, y_N)\}$, $t = \infty$ (min distance)
- **Input:** a test data input x
- **KNN classification rule (implementation):**
 - Denote by $\text{KNN}(x)$ the set of the K nearest neighbors of x in D_{tr}
 - $\text{KNN}(x) \subset D_{tr}$, $|\text{KNN}(x)| = K$
 - $\forall (x', y') \in D_{tr} \setminus \text{KNN}(x)$, $\text{dist}(x, x') \geq \max_{(x'', y'') \in \text{KNN}(x)} \text{dist}(x, x'')$
 - Every point in D_{tr} but not in $\text{KNN}(x)$ is at least as far away from x as the furthest point in $\text{KNN}(x)$
 - **Return:** $\hat{y} = \hat{h}(x) = \text{mode}(\{y''; (x'', y'') \in \text{KNN}(x)\})$: most frequent label in $\text{KNN}(x)$

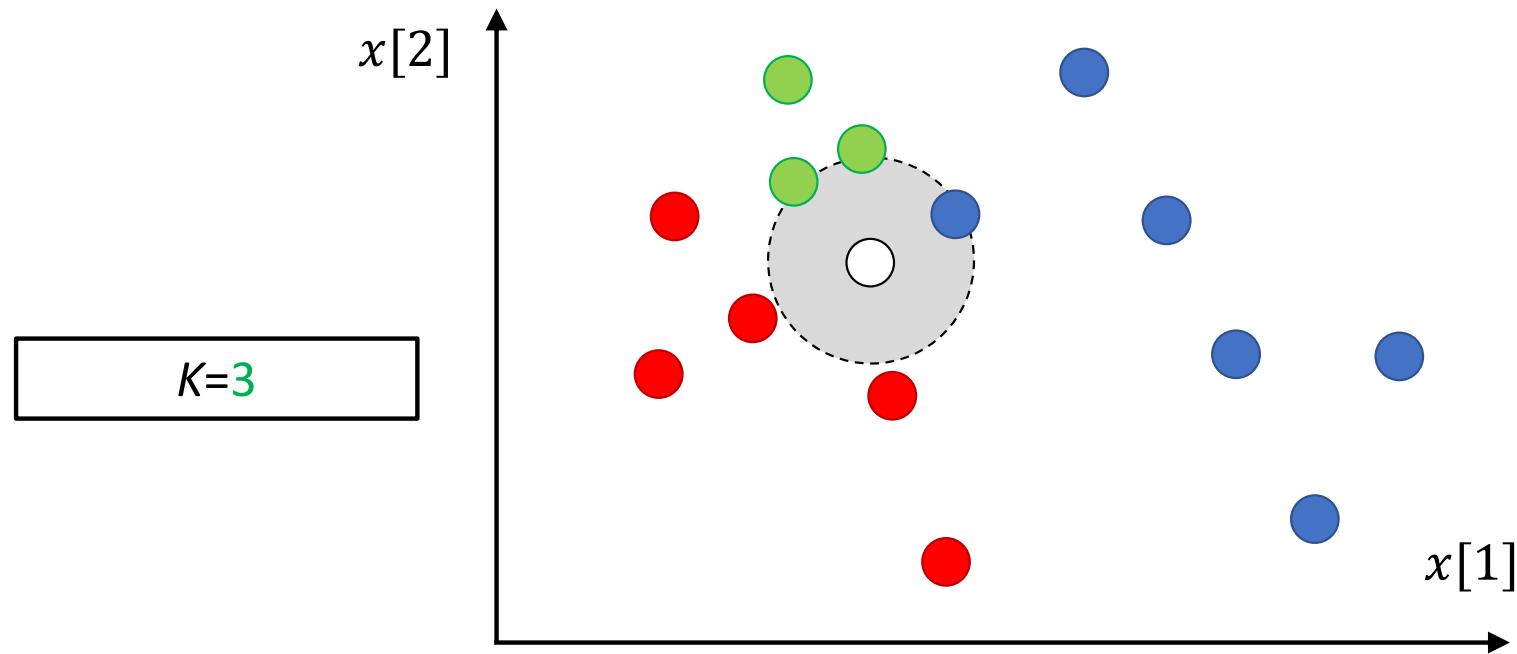
Question

- How to resolve the case of a tie?
- Applicable to binary classification only or not?

KNN (Euclidean distance) for multiple classes



KNN (Euclidean distance) for multiple classes



Variants of KNN: Rules

Weighted KNN

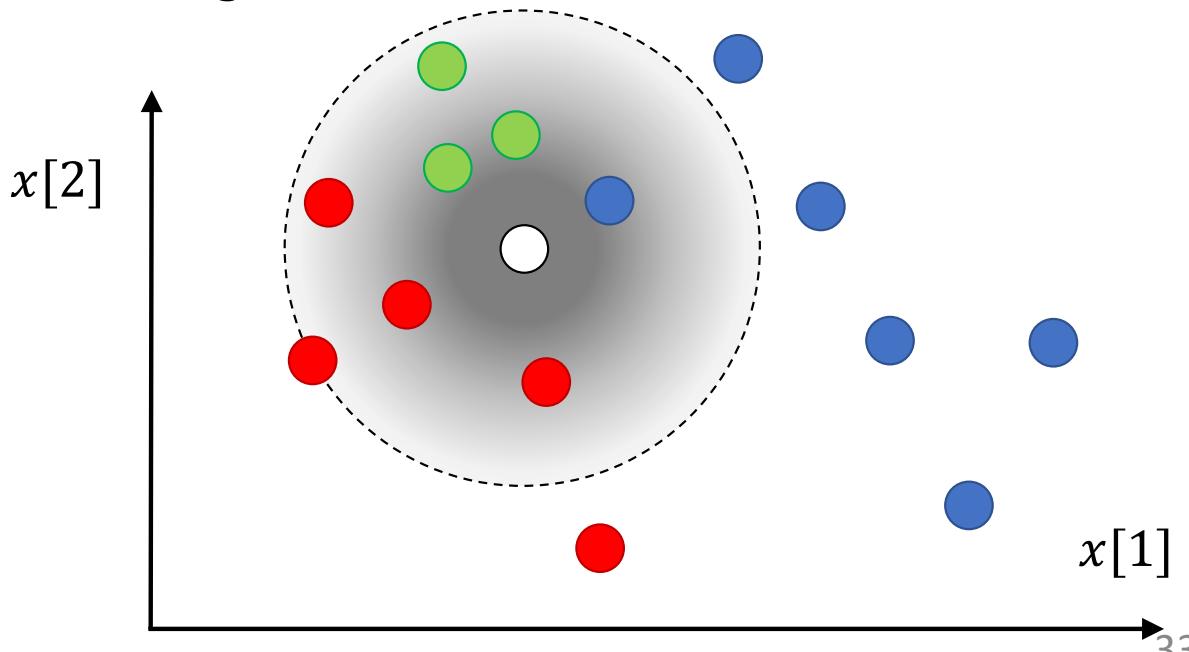
- for class $c = 1:C$
 - $\text{Vote}(c) = \sum_{(x',y') \in KNN(x)} \mathbf{1}[y' == c] \times \text{weight}(x, x')$

- $\hat{y} = \underset{c}{\operatorname{argmax}} \text{Vote}(c)$

- For example

- $\text{weight}(x, x') = e^{-\frac{\|x-x'\|^2}{2\sigma^2}}$

- Radial basis function (RBF) kernel



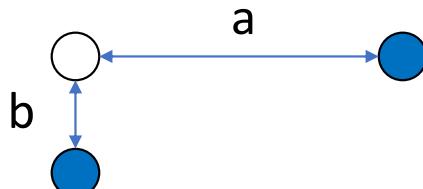
Variants of KNN: Distances

- Cosine distance
 - $\cos(\mathbf{x}, \mathbf{x}') = \frac{\mathbf{x}^T \mathbf{x}'}{\|\mathbf{x}\|_2 \|\mathbf{x}'\|_2}$
 - $1 - \cos(\mathbf{x}, \mathbf{x}') = 1 - \frac{\mathbf{x}^T \mathbf{x}'}{\|\mathbf{x}\|_2 \|\mathbf{x}'\|_2}$
- Mahalanobis distance (less strict form)
 - $(\mathbf{x} - \mathbf{x}')^T \mathbf{A} (\mathbf{x} - \mathbf{x}')$, where $\mathbf{A} \in \mathbb{R}^{D \times D}$ is symmetric positive semi-definite (PSD)
 - PSD: $\forall \mathbf{x} \neq \mathbf{0}, \mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0$
 - $(\mathbf{x} - \mathbf{x}')^T \mathbf{A} (\mathbf{x} - \mathbf{x}') = (\mathbf{x} - \mathbf{x}')^T \mathbf{M}^T \mathbf{M} (\mathbf{x} - \mathbf{x}') = (\mathbf{M}(\mathbf{x} - \mathbf{x}'))^T \mathbf{M}(\mathbf{x} - \mathbf{x}')$
 - $\mathbf{M} \in \mathbb{R}^{D' \times D}$ and has no other constraint

Variants of KNN: Distances

- What does Mahalanobis distance mean? How to set M ?

$$M = \begin{bmatrix} 1 & 0 \\ \frac{1}{a} & 0 \\ 0 & 1 \\ 0 & \frac{1}{b} \end{bmatrix}$$



$$\begin{bmatrix} a \\ 0 \\ 0 \\ b \end{bmatrix}^T \begin{bmatrix} 1 & 0 \\ \frac{1}{a} & 0 \\ 0 & 1 \\ 0 & \frac{1}{b} \end{bmatrix}^T \begin{bmatrix} 1 & 0 \\ \frac{1}{a} & 0 \\ 0 & 1 \\ 0 & \frac{1}{b} \end{bmatrix} \begin{bmatrix} a \\ 0 \\ 0 \\ b \end{bmatrix} = 1,$$

$$\begin{bmatrix} 0 \\ b \\ 0 \\ 0 \end{bmatrix}^T \begin{bmatrix} 1 & 0 \\ \frac{1}{a} & 0 \\ 0 & 1 \\ 0 & \frac{1}{b} \end{bmatrix}^T \begin{bmatrix} 1 & 0 \\ \frac{1}{a} & 0 \\ 0 & 1 \\ 0 & \frac{1}{b} \end{bmatrix} \begin{bmatrix} 0 \\ b \\ 0 \\ 0 \end{bmatrix} = 1$$

- Long-standing statistical problem. Can estimate M or A from covariance
- Distance metric learning to learn M or A

Commonly-used tricks for pre-processing

- Z-score
 - Compute training mean: $\mu = \frac{1}{N} \sum_{n=1}^N x_n$
 - Compute training STD per dimension: $\rho[d] = \sqrt{\frac{1}{N-1} \sum_{n=1}^N (x_n[d] - \mu)^2}$
 - For every data instances x , training or test, do $x[d] \leftarrow \frac{x[d] - \mu[d]}{\rho[d]}$
 - Make the training data “zero”-mean, and equal variance per dimension
- In practice, many other pre-processing methods like L2-normalization

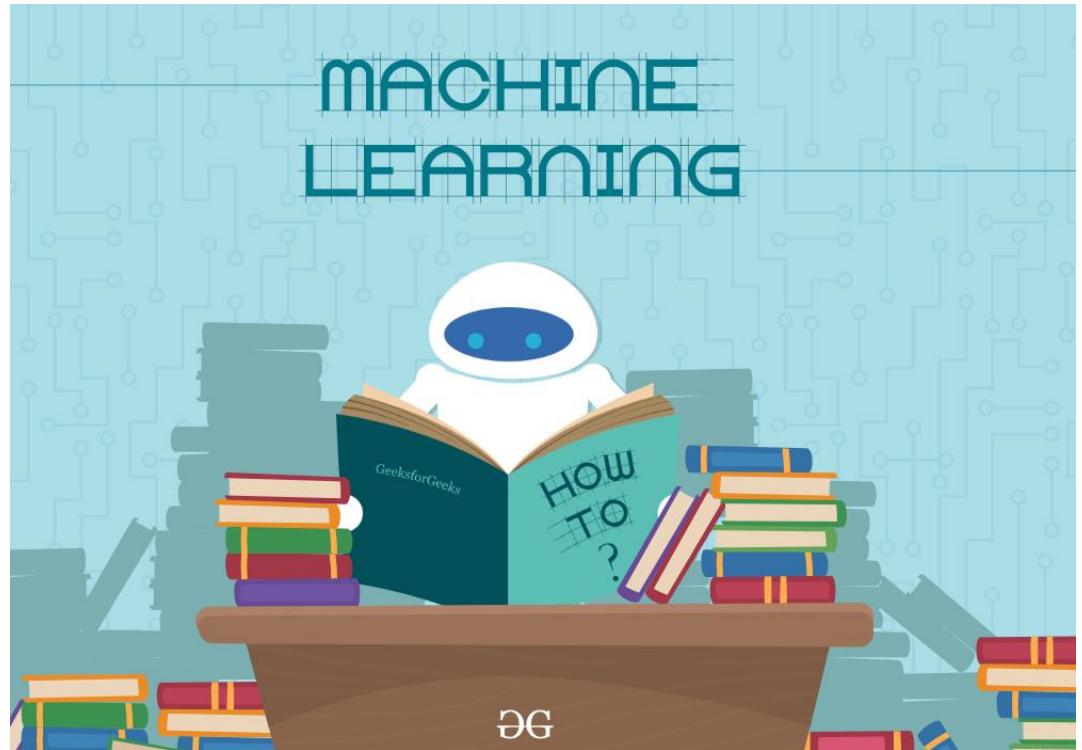
Today

Review

- ML problem solving
- Training, validation, testing

Nearest neighbor

- Basic and variants
- Properties
- Curse of dimensionality
- Bayes optimal classifier



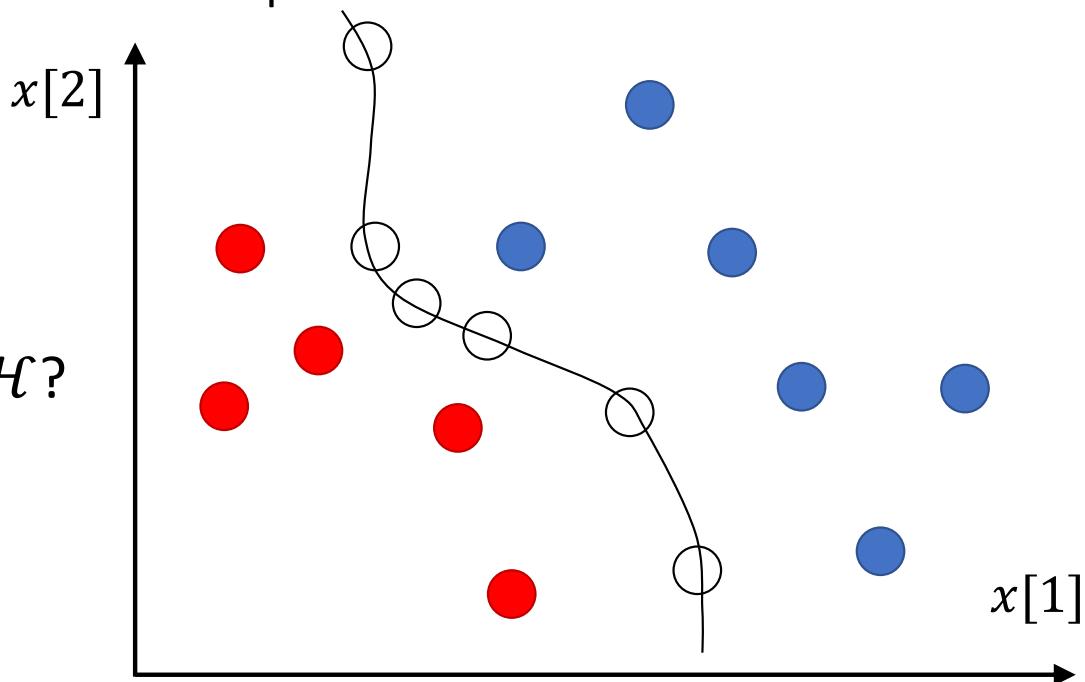
Properties

- Does learning EVER happen?
 - Seems NO before seeing the test data instance \mathbf{x}
 - **Lazy learning:** generalization of the training data is delayed until a test data comes.
 - Opposite: eager learning, which generalizes the training data before test data.
- Is KNN merely memorization?
 - Recall: $\hat{h}(\mathbf{x}) = \begin{cases} y_n, & \text{if } \exists (\mathbf{x}_n, y_n) \in D_{tr} \text{ s.t. } \mathbf{x} = \mathbf{x}_n \\ 0, & \text{otherwise} \end{cases}$
 - NN: Recall: $\hat{h}(\mathbf{x}) = \begin{cases} y_n, & \text{if } \exists (\mathbf{x}_n, y_n) \in D_{tr} \text{ s.t. } \mathbf{x} = \mathbf{x}_n \\ y_m & \text{s.t. } m = \operatorname{argmin}_n \operatorname{dist}(\mathbf{x}, \mathbf{x}_n) \end{cases}$
 - KNN: Recall: $\hat{h}(\mathbf{x}) = \operatorname{mode}(\{y''; (\mathbf{x}'', y'') \in \operatorname{KNN}(\mathbf{x})\})$

Assuming no
duplicated
data input!

Properties

- Why does KNN not just memorize?
 - Similar Inputs (feature vectors) have similar outputs (labels)
 - Every successful ML algorithms MUST make assumptions
- What is the hypothesis class \mathcal{H} ?
 - Linear?
- What are the controlling factors of \mathcal{H} ?
 - K
 - Distance metric
 - They are called “hyper-parameters”



Properties

- What is the training error of NN?

- $\hat{\epsilon}(\hat{h}; D_{tr}) = \frac{1}{N} \sum_{n=1}^N \ell(y_n, \hat{h}(\mathbf{x}_n; D_{tr}))$

- Assume there are no pair of training data share the same input

- $\hat{\epsilon}(\hat{h}; D_{tr}) = \frac{1}{N} \sum_{n=1}^N \ell(y_n, \hat{y}_n) = 0$

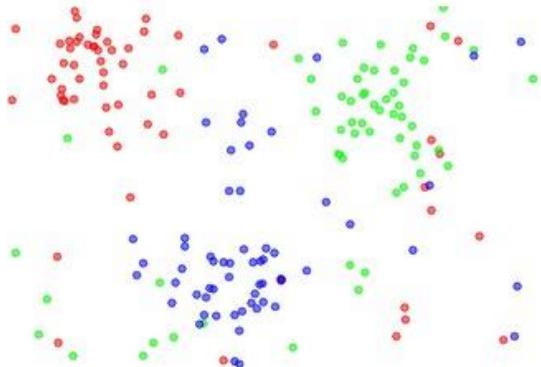
- What is the training error of KNN?

- May NOT be 0

What is the effect of K ?

What is the effect of K ?

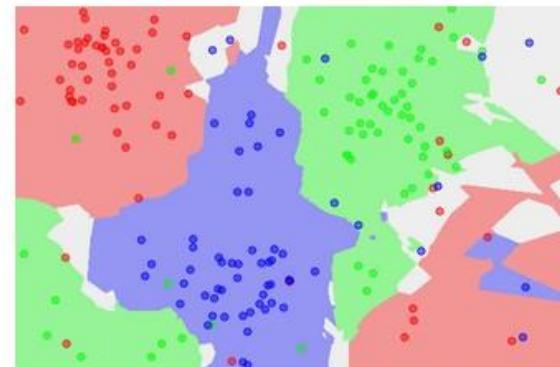
the data



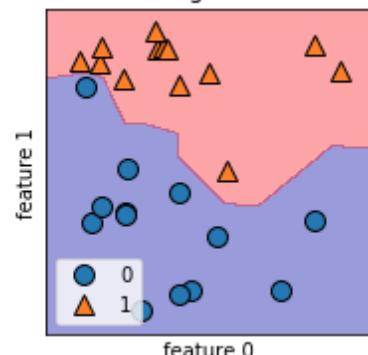
NN classifier



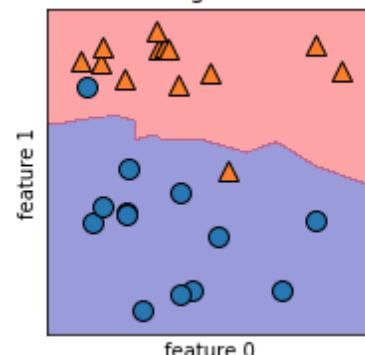
5-NN classifier



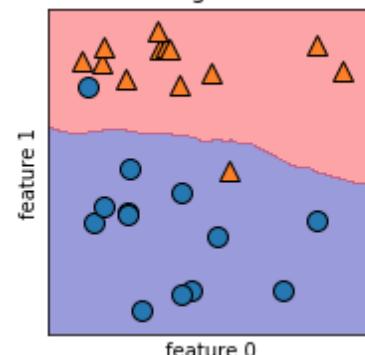
1 neighbor(s)



3 neighbor(s)



9 neighbor(s)

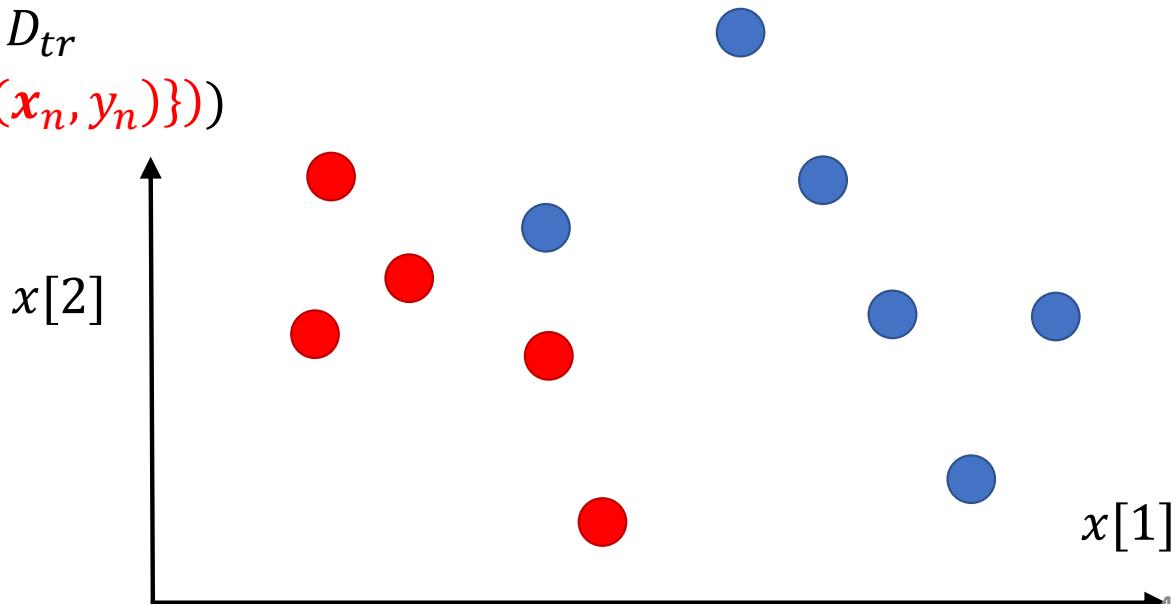


How to choose K ?

Leave-one-out Cross-validation

Each training instance takes turns to temporally be the “single” Val instance

- $\hat{\epsilon}(\hat{h}; D_{val}) = 0$
- for $n = 1 : N$
 - temporally remove (x_n, y_n) from D_{tr}
 - $\hat{\epsilon}(\hat{h}; D_{val}) += \ell(y_n, \hat{h}(x_n; D_{tr} \setminus \{(x_n, y_n)\}))$
- $\hat{\epsilon}(\hat{h}; D_{val}) = \frac{\hat{\epsilon}(\hat{h}; D_{val})}{N}$



Pros and Cons

- Pros

- Easy to implement (and have theoretical guarantee)
- No training, suitable for “dynamically” changing training data
- Suitable for “unknown” number of classes
- Suitable for not well-defined classes (e.g., image retrieval, Google image search)



Okapi un animal
misterios Okapia
johnstoni - zooland.ro

- Cons

- Need to carry all the training data to testing
- $O(NxD)$ test time (speed-up by KD-tree or hashing)
- Become less effective in high-dimensional space: Curse of dimensionality
- **Popular distance metrics** may not be good enough

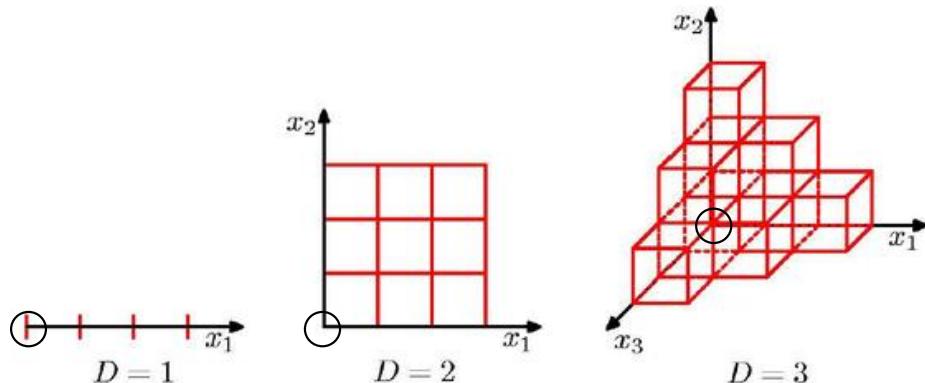
Questions?



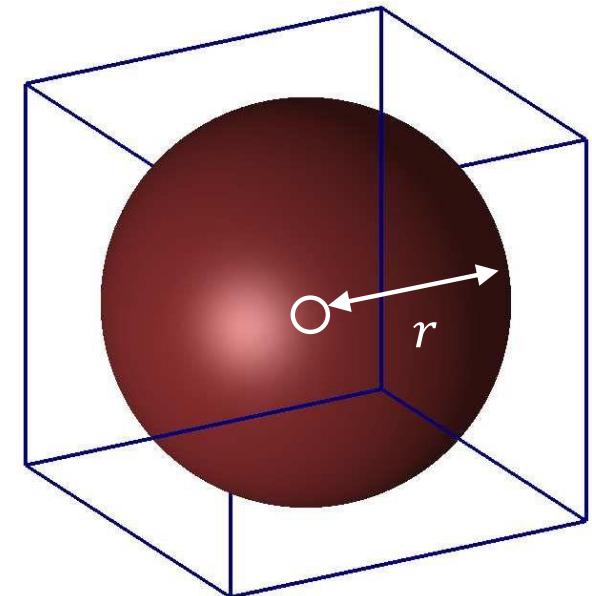
THE OHIO STATE UNIVERSITY

Curse of dimensionality

- In high-dimensional space
 - Hard to find training data that are close to the test data



Higher dimension, hard to find neighbors!

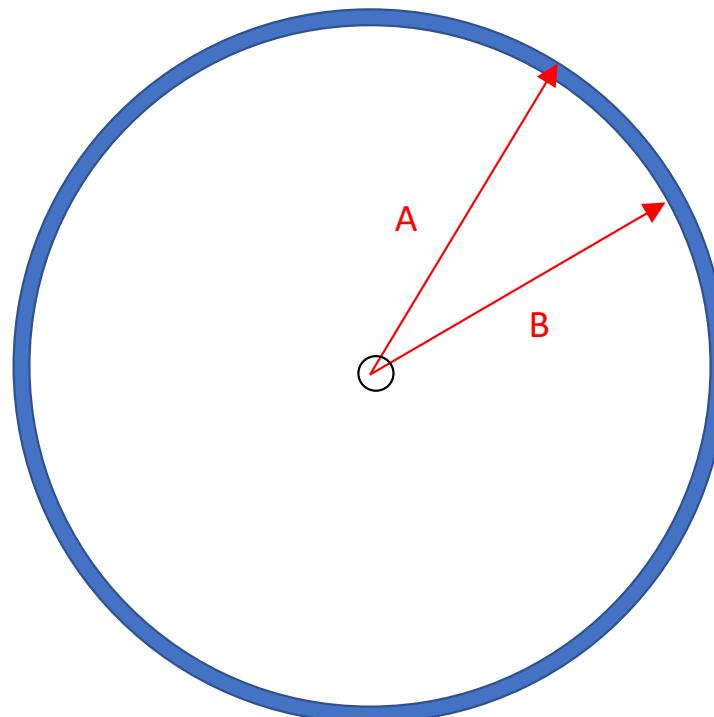


Most uniformly distributed data are outside the circle

$$V(D) = \frac{\pi^{\frac{D}{2}}}{\Gamma(\frac{D}{2} + 1)} r^D \quad \text{v.s.} \quad (2r)^D$$

Curse of dimensionality

- In high-dimensional space
 - Every training data seem to be equally far away from the test data



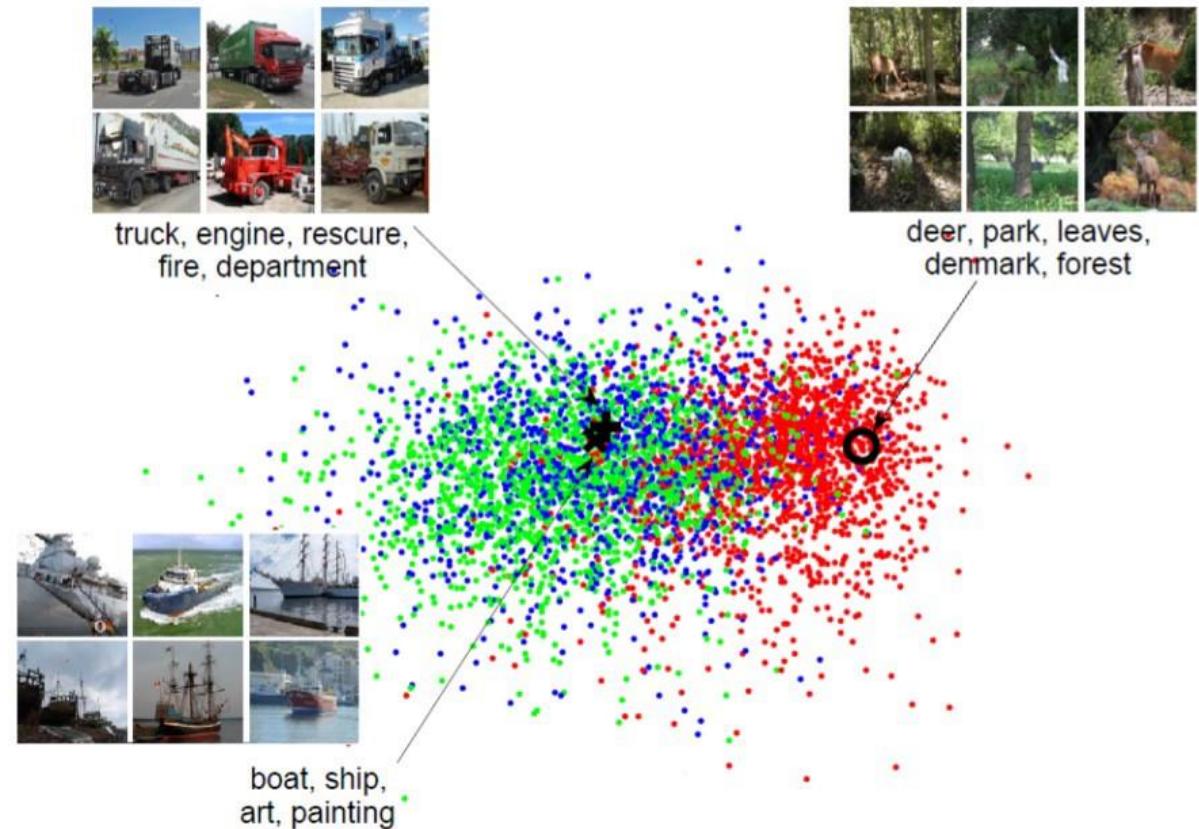
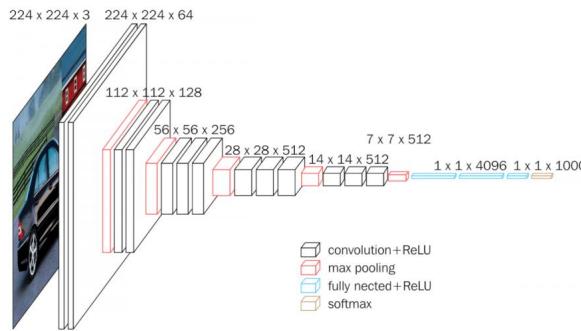
$$V(A) = \frac{\pi^{\frac{D}{2}}}{\Gamma(\frac{D}{2} + 1)} (r + \varepsilon)^D$$

$$V(B) = \frac{\pi^{\frac{D}{2}}}{\Gamma(\frac{D}{2} + 1)} r^D$$

$$\frac{V(A) - V(B)}{V(B)} \uparrow, \text{ as } D \uparrow$$

Feature learning, dimensionality reduction

$$x = g(I)$$
$$g: \mathcal{I} \mapsto \mathcal{X}$$



Questions?



THE OHIO STATE UNIVERSITY

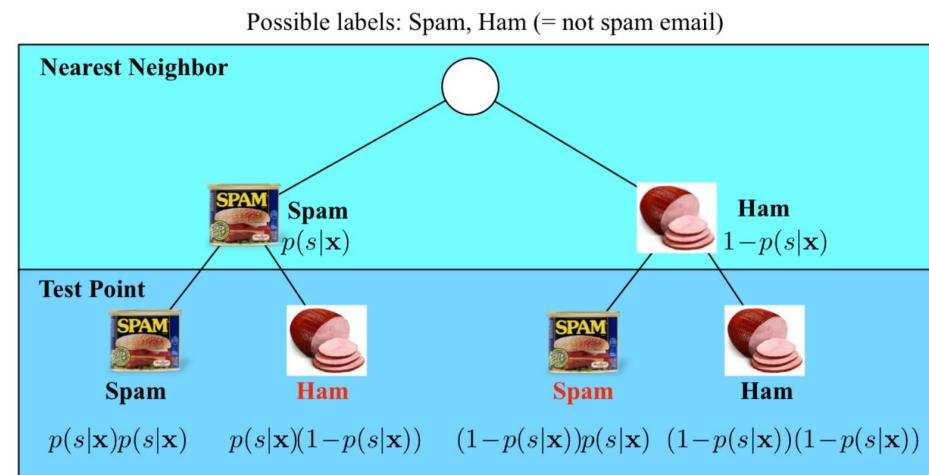
Theoretical guarantee: Loss lower bound

- Assume our true data distribution is $P(X \in \mathcal{X}, Y \in \mathcal{Y}) = P(X)P(Y|X)$
- Assume you know the distribution (impossible mostly, thus an ideal case):
 - What will you predict given \mathbf{x} ? Consider the 0-1 loss!
 - Predict $y^* = \underset{y}{\operatorname{argmax}} P(Y = y|\mathbf{x})$, what is the “expected” loss?
 - $(1 - P(Y = y^*|\mathbf{x}))$
 - Predicting any other label will lead to a higher expected loss
- **Bayes optimal classifier:** $y^* = \underset{y}{\operatorname{argmax}} P(Y = y|\mathbf{x})$
 - “Generalization” loss: $E_{\mathbf{x} \sim P(\mathbf{x})}[1 - P(Y = y^*|\mathbf{x})]$
 - Cannot go even lower without other extra information

Theoretical guarantee: NN vs. Bayes optimal

- Claim: As $N \rightarrow \infty$, the NN classification error is no more than **TWICE** of Bayes optimal

- When $N \rightarrow \infty$, the nearest neighbor x' of x has $\text{dist}(x, x') \rightarrow 0$, which means x' is almost identical to x
- Assign the label of x' to x
- When will it go wrong? (let's try binary)
 - $x': +1, x: -1$
 - $x': -1, x: +1$
 - $2 \times P(Y = y^*|x)(1 - P(Y = y^*|x))$
 - $\leq 2 \times (1 - P(Y = y^*|x))$



- $\text{Loss}(\text{KNN}) \leq 2 \times (1 - P(Y = y^*|x)) = 2 \times \text{Loss}(\text{Bayes optimal})$

Summary

- KNN classifiers are lazy learning methods
 - The same algorithm works for binary and multi-class classification
 - Assumptions: similar inputs, similar outputs
-
- Different hypothesis classes are governed by (1) K, (2) distance metric
 - “Leave-One-Out” (LOO) cross-validation
-
- Pros: Easy to implement, “no training”, **theoretical guarantee**
 - Cons: test time, (training) data memory, **curse of dimensionality**

CSE 5523: Nearest Neighbor + Linear Regression



THE OHIO STATE UNIVERSITY

Important for this week

- **Register** for the class on Piazza: our platform for discussion and communication
 - https://piazza.com/osu/autumn2024/cse5523_chao (please use name.#@osu.edu)
 - Access code: osu-cse-5523-AU24-chao
- **Math review/self-diagnostic:** do Homework #0 and suggested materials on the website --- **extremely important to check your readiness for the course**
- **Office hours:** start this Thursday
 - Mine: Tuesday 11 am - 12 pm, Thursday 6 pm - 7 pm (DL 587)
 - TA (Tai-Yu Pan): Monday 11 am - 12 pm, Friday 4 pm - 5 pm (BE406, #12 and #13)

Homework

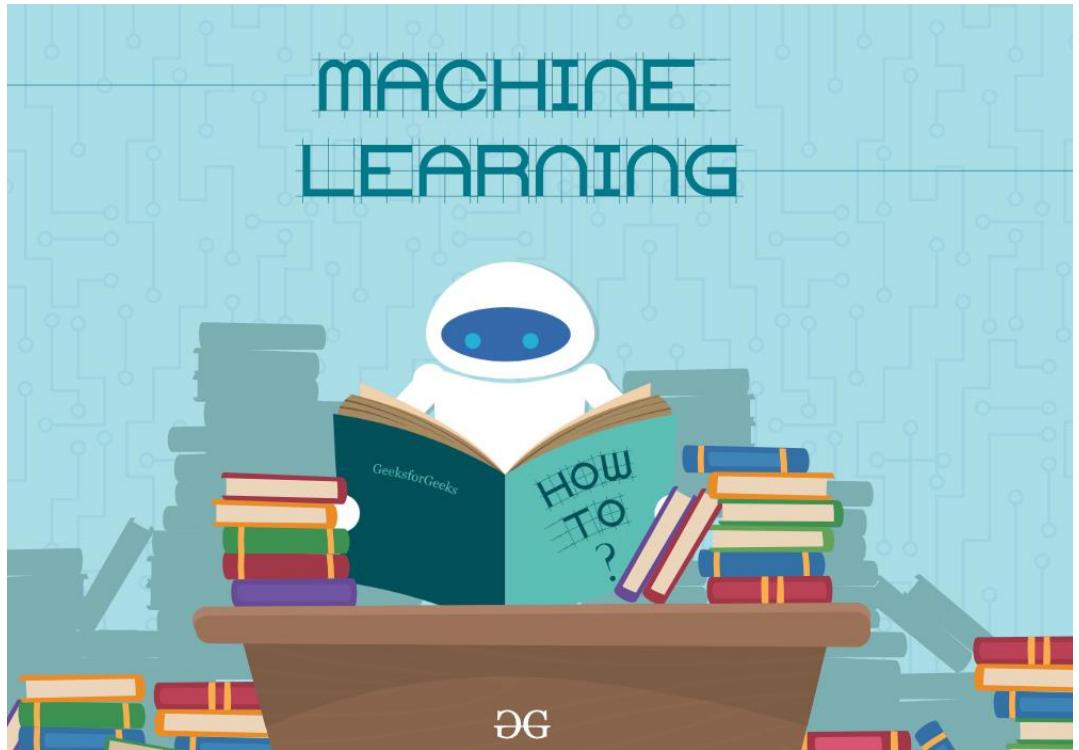
- HW # 1:
 - Release date: 9/3
 - Due date: 9/17 (11:59 pm)
 - Problem set + Programming set

Today

Nearest neighbor

- Review
- Properties (continued)
- Curse of dimensionality
- Bayes optimal classifier

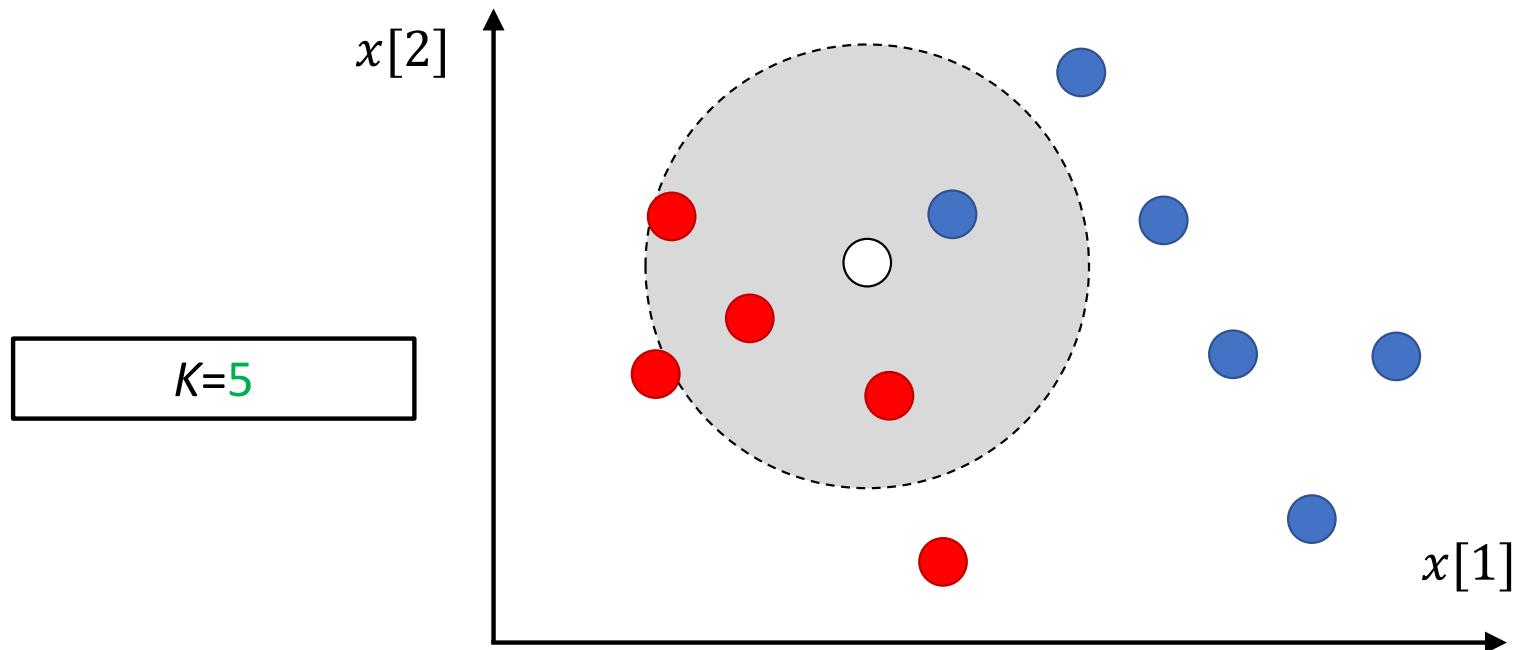
Linear regression



K-Nearest neighbors (KNN)

- What is the assumption?
 - Similar Inputs (feature vectors) have similar outputs (labels)
 - This makes KNN not merely memorization!
- Classification rule:
 - For a test data input x , assign the most common label amongst its K most similar training data instances in D_{tr}
 - aka, majority vote

KNN (Euclidean distance)



$$\|x - x_n\|_2 = \left(\sum_{d=1}^D |x[d] - x_n[d]|^2 \right)^{1/2}$$

Distance

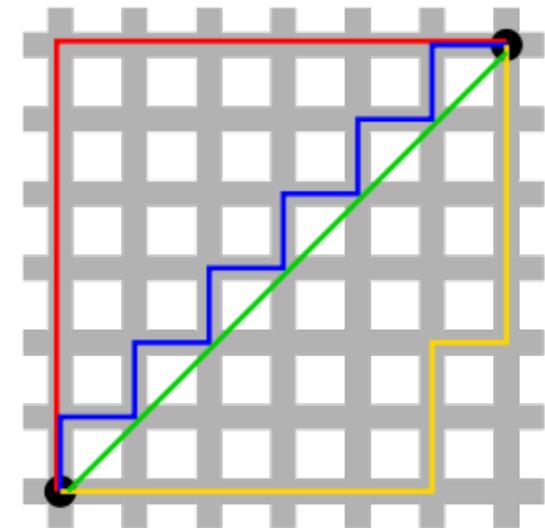
- The better the distance metric reflects label similarity, x : 2-dimensional vectors the better the classifier is.

- L_2 distance (Euclidean): $\|x - x_n\|_2 = \left(\sum_{d=1}^D |x[d] - x_n[d]|^2 \right)^{1/2}$

- L_1 distance: $\|x - x_n\|_1 = \sum_{d=1}^D |x[d] - x_n[d]|$

- L_p norm: $\|x - x_n\|_p = \left(\sum_{d=1}^D |x[d] - x_n[d]|^p \right)^{1/p}$

- Cosine distance, Mahalanobis distance,
- Data pre-processing (e.g., z-score)



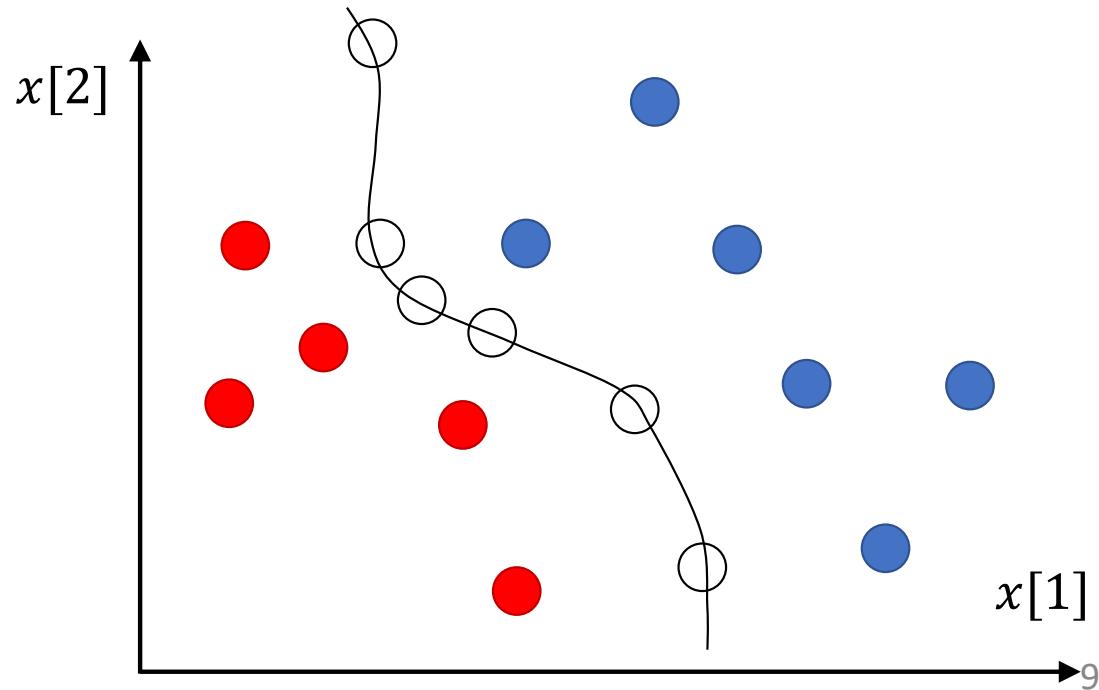
Green line is Euclidean distance.
Red, Blue, and Yellow lines are L_1 distance

Commonly-used tricks for pre-processing

- Z-score
 - Compute training mean: $\mu = \frac{1}{N} \sum_{n=1}^N x_n$
 - Compute training STD per dimension: $\rho[d] = \sqrt{\frac{1}{N-1} \sum_{n=1}^N (x_n[d] - \mu)^2}$
 - For every data instances x , training or test, do $x[d] \leftarrow \frac{x[d] - \mu[d]}{\rho[d]}$
 - Make the training data “zero”-mean, and equal variance per dimension
- In practice, many other pre-processing methods like L2-normalization

Properties

- Lazy learning methods
- Work for binary and multi-class classification
- Different hypothesis classes:
 - K + distance metric
 - They are called “hyper-parameters”

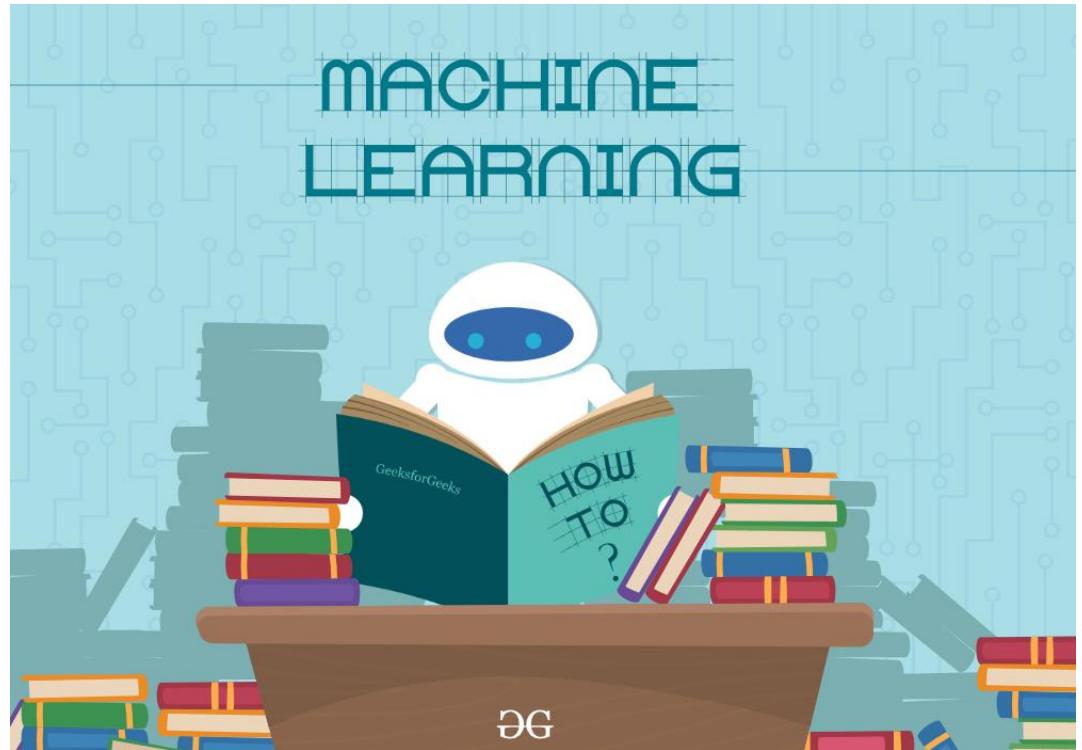


Today

Nearest neighbor

- Review
- Properties (continued)
- Curse of dimensionality
- Bayes optimal classifier

Linear regression



Properties

- What is the training error of NN?

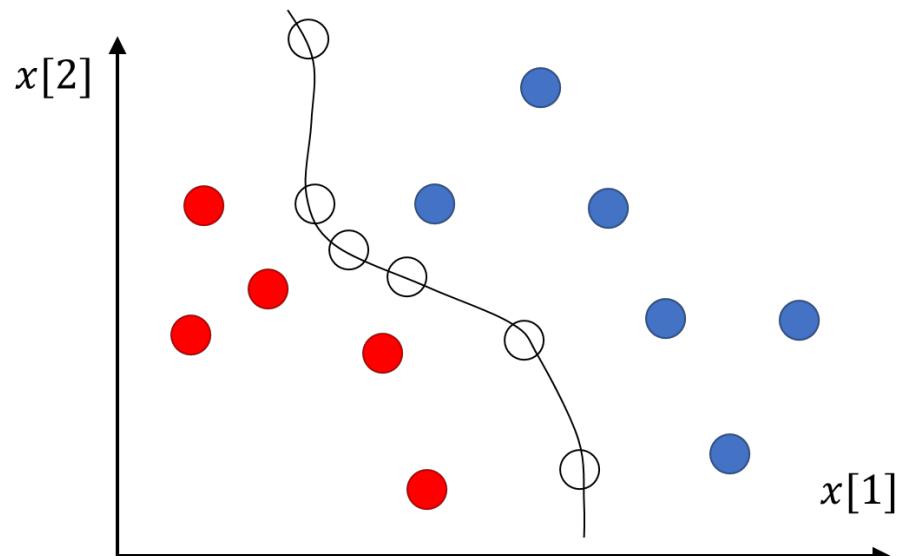
- $\hat{\epsilon}(\hat{h}; D_{tr}) = \frac{1}{N} \sum_{n=1}^N \ell(y_n, \hat{h}(x_n; D_{tr}))$

- Assume there are no pair of training data share the same input

- $\hat{\epsilon}(\hat{h}; D_{tr}) = \frac{1}{N} \sum_{n=1}^N \ell(y_n, \hat{y}_n) = 0$

- What is the training error of KNN?

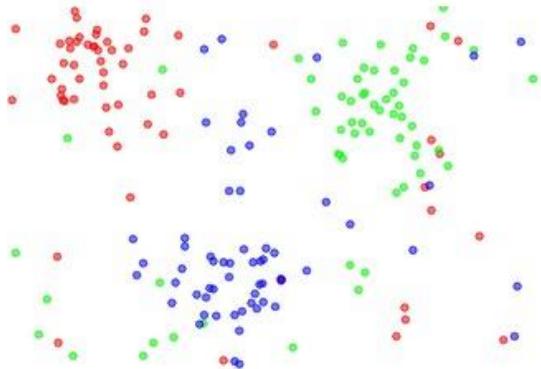
- May NOT be 0



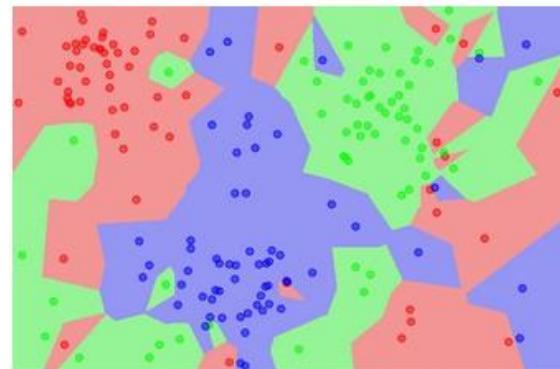
What is the effect of K ?

What is the effect of K ?

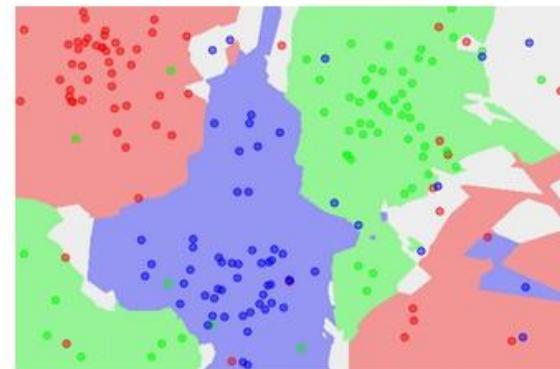
the data



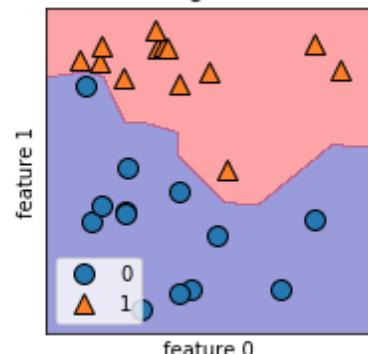
NN classifier



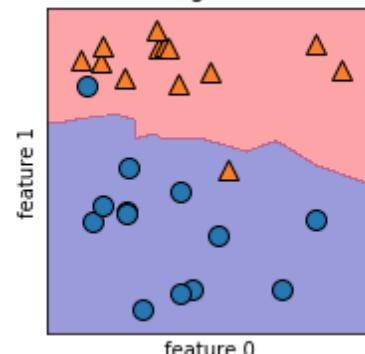
5-NN classifier



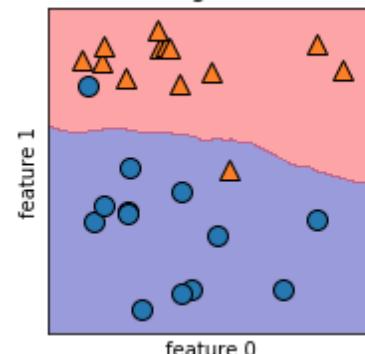
1 neighbor(s)



3 neighbor(s)



9 neighbor(s)

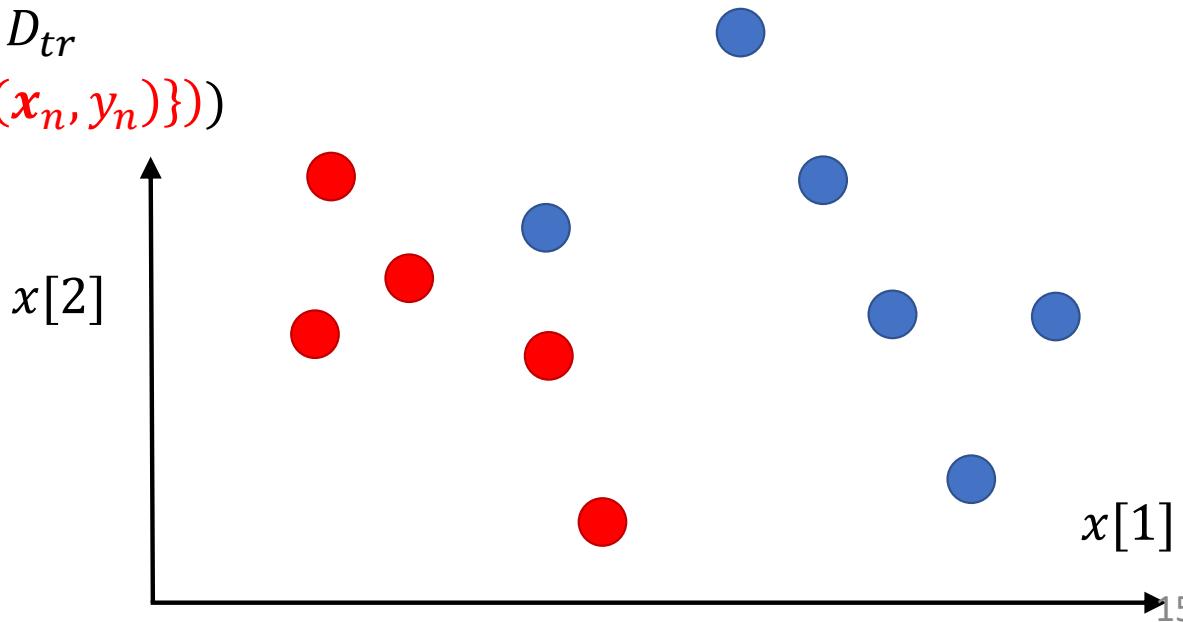


How to choose K ?

Leave-one-out Cross-validation

Each training instance takes turns to temporally be the “single” Val instance

- $\hat{\epsilon}(\hat{h}; D_{val}) = 0$
- for $n = 1 : N$
 - temporally remove (x_n, y_n) from D_{tr}
 - $\hat{\epsilon}(\hat{h}; D_{val}) += \ell(y_n, \hat{h}(x_n; D_{tr} \setminus \{(x_n, y_n)\}))$
- $\hat{\epsilon}(\hat{h}; D_{val}) = \frac{\hat{\epsilon}(\hat{h}; D_{val})}{N}$



Pros and Cons

- Pros

- Easy to implement (and have theoretical guarantee)
- No training, suitable for “dynamically” changing training data
- Suitable for “unknown” number of classes
- Suitable for not well-defined classes (e.g., image retrieval, Google image search)



Okapi un animal
misterios Okapia
johnstoni - zooland.ro

- Cons

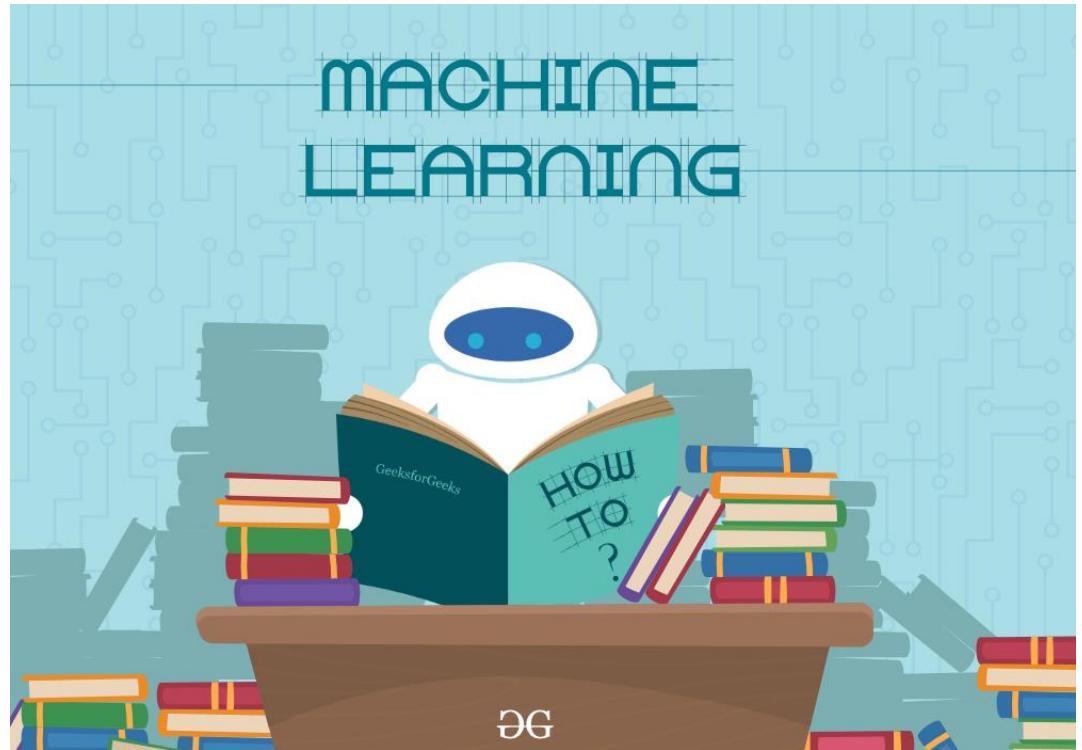
- Need to carry all the training data to testing
- $O(NxD)$ test time (speed-up by KD-tree or hashing)
- Become less effective in high-dimensional space: Curse of dimensionality
- **Popular distance metrics** may not be good enough

Today

Nearest neighbor

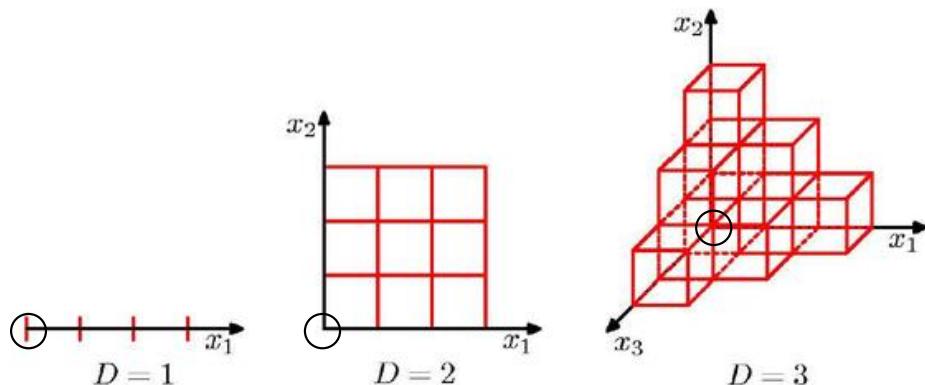
- Review
- Properties (continued)
- Curse of dimensionality
- Bayes optimal classifier

Linear regression

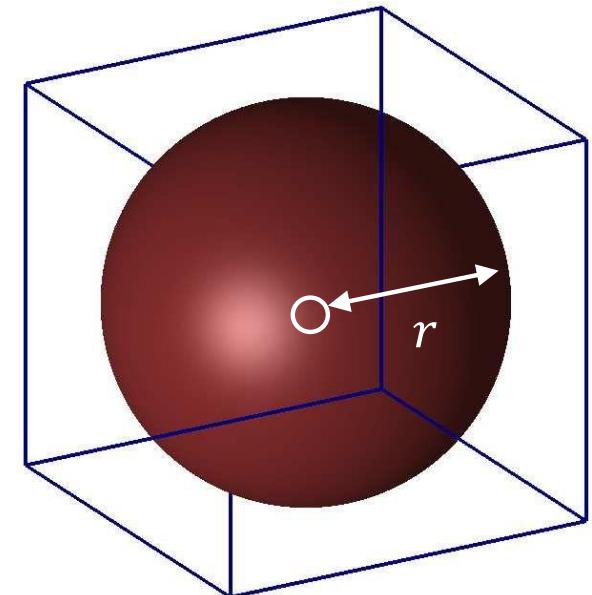


Curse of dimensionality

- In high-dimensional space
 - Hard to find training data that are close to the test data



Higher dimension, hard to find neighbors!

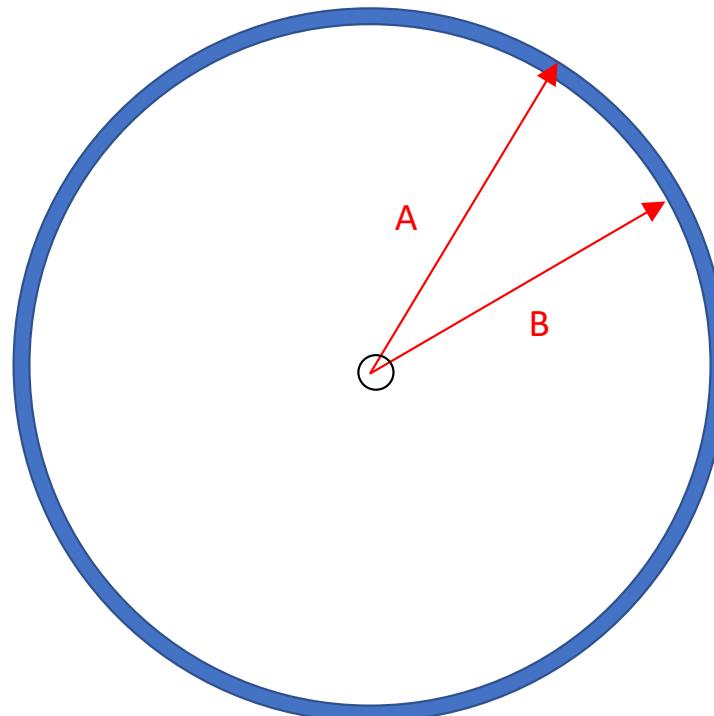


Most uniformly distributed data are outside the circle

$$V(D) = \frac{\pi^{\frac{D}{2}}}{\Gamma(\frac{D}{2} + 1)} r^D \quad \text{v.s.} \quad (2r)^D$$

Curse of dimensionality

- In high-dimensional space
 - Every training data seem to be equally far away from the test data



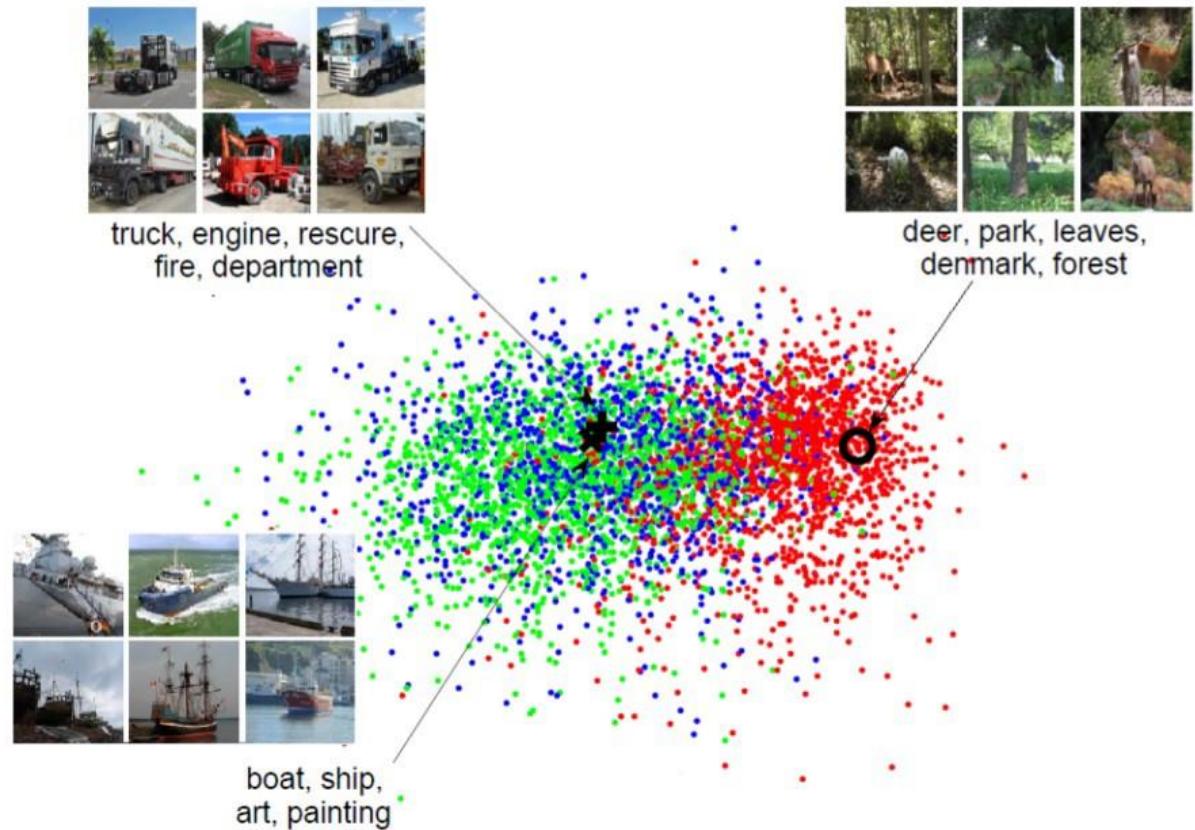
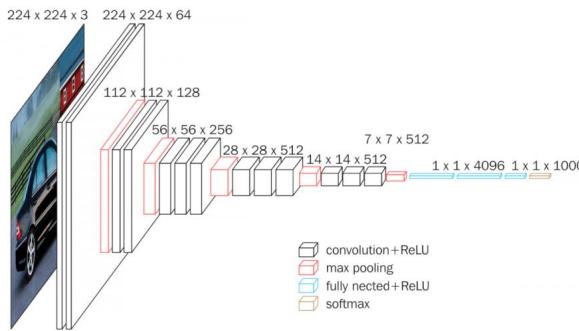
$$V(A) = \frac{\pi^{\frac{D}{2}}}{\Gamma(\frac{D}{2} + 1)} (r + \varepsilon)^D$$

$$V(B) = \frac{\pi^{\frac{D}{2}}}{\Gamma(\frac{D}{2} + 1)} r^D$$

$$\frac{V(A) - V(B)}{V(B)} \uparrow, \text{ as } D \uparrow$$

Feature learning, dimensionality reduction

$$x = g(I)$$
$$g: \mathcal{I} \mapsto \mathcal{X}$$

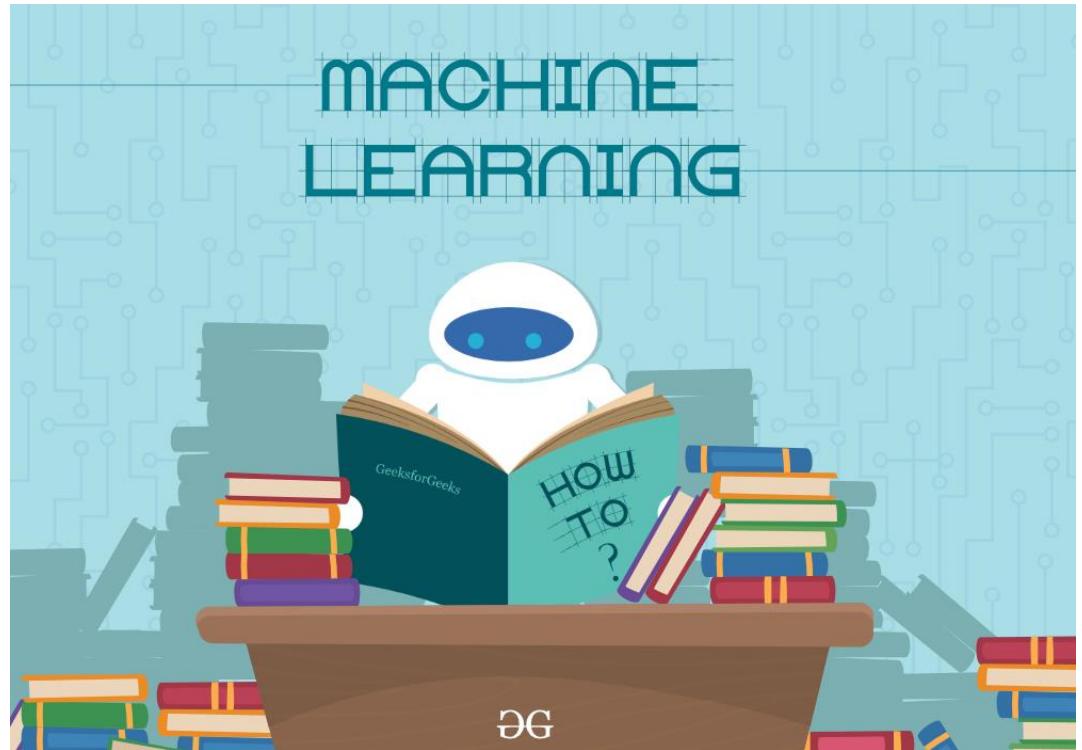


Today

Nearest neighbor

- Review
- Properties (continued)
- Curse of dimensionality
- Bayes optimal classifier

Linear regression



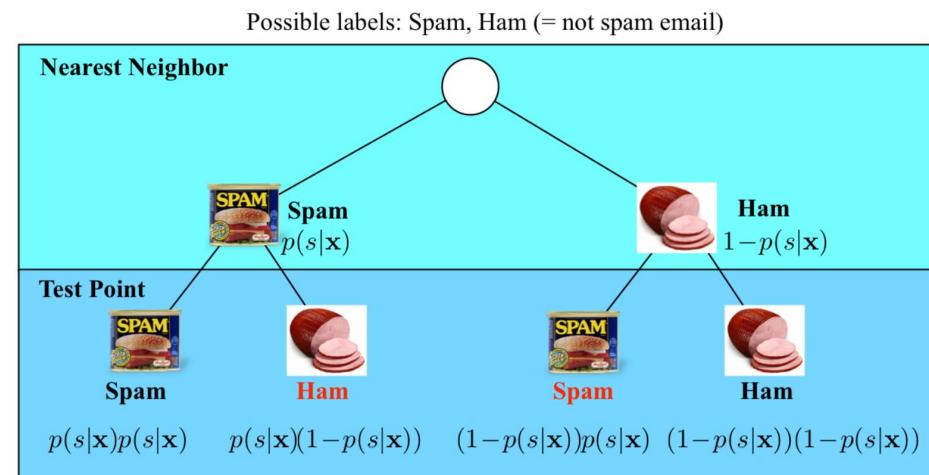
Theoretical guarantee: Loss lower bound

- Assume our true data distribution is $P(X \in \mathcal{X}, Y \in \mathcal{Y}) = P(X)P(Y|X)$
- Assume you know the distribution (impossible mostly, thus an ideal case):
 - What will you predict given \mathbf{x} ? Consider the 0-1 loss!
 - Predict $y^* = \underset{y}{\operatorname{argmax}} P(Y = y|\mathbf{x})$, what is the “expected” loss?
 - $(1 - P(Y = y^*|\mathbf{x}))$
 - Predicting any other label will lead to a higher expected loss
- **Bayes optimal classifier:** $y^* = \underset{y}{\operatorname{argmax}} P(Y = y|\mathbf{x})$
 - “Generalization” loss: $E_{\mathbf{x} \sim P(\mathbf{x})}[1 - P(Y = y^*|\mathbf{x})]$
 - Cannot go even lower without other extra information

Theoretical guarantee: NN vs. Bayes optimal

- Claim: As $N \rightarrow \infty$, the NN classification error is no more than **TWICE** of Bayes optimal

- When $N \rightarrow \infty$, the nearest neighbor x' of x has $\text{dist}(x, x') \rightarrow 0$, which means x' is almost identical to x
- Assign the label of x' to x
- When will it go wrong? (let's try binary)
 - $x': +1, x: -1$
 - $x': -1, x: +1$
 - $2 \times P(Y = y^*|x)(1 - P(Y = y^*|x))$
 - $\leq 2 \times (1 - P(Y = y^*|x))$



- $\text{Loss}(\text{KNN}) \leq 2 \times (1 - P(Y = y^*|x)) = 2 \times \text{Loss}(\text{Bayes optimal})$

Summary

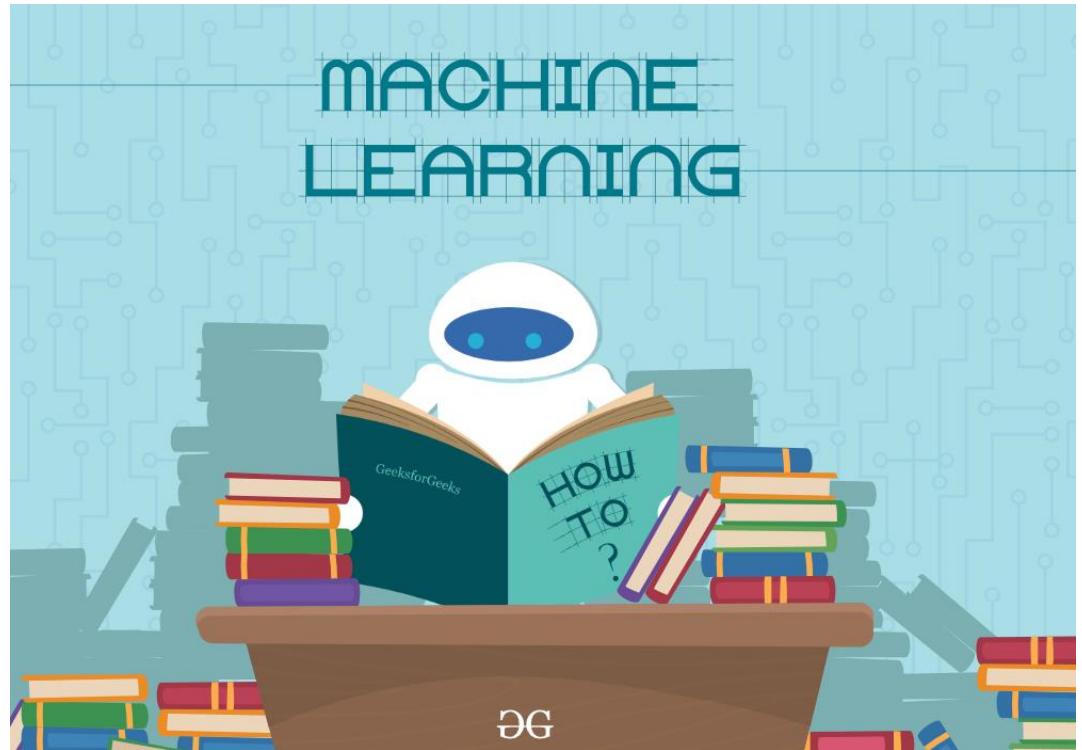
- KNN classifiers are lazy learning methods
 - The same algorithm works for binary and multi-class classification
 - Assumptions: similar inputs, similar outputs
-
- Different hypothesis classes are governed by (1) K, (2) distance metric
 - “Leave-One-Out” (LOO) cross-validation
-
- Pros: Easy to implement, “no training”, **theoretical guarantee**
 - Cons: test time, (training) data memory, **curse of dimensionality**

Today

Nearest neighbor

- Review
- Properties (continued)
- Curse of dimensionality
- Bayes optimal classifier

Linear regression



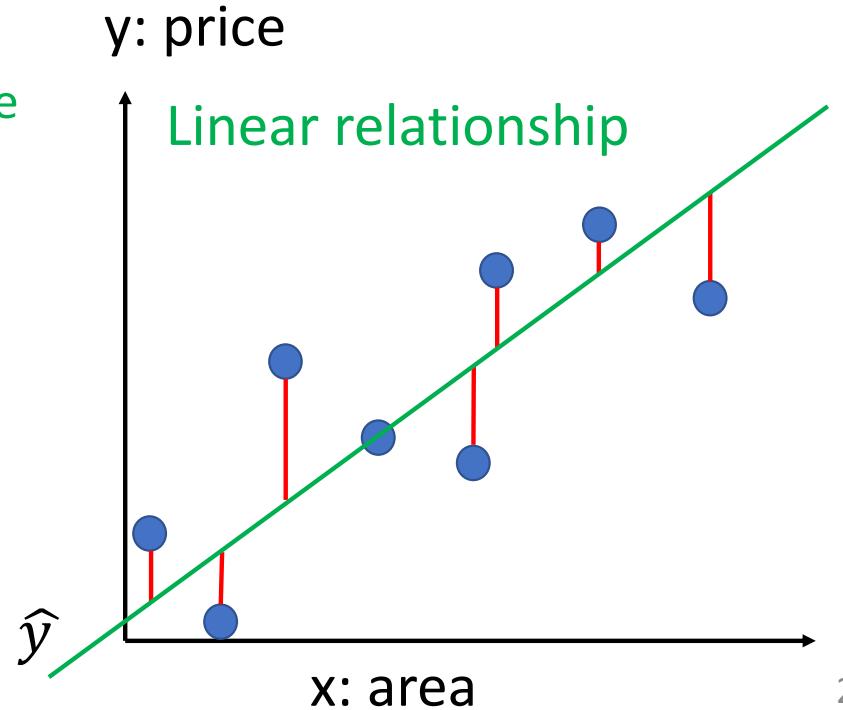
(Linear) regression

- Predicting a continuous output
 - Price, distance, location, ...
- Setup
 - $D_{tr} = \{(x_n \in \mathbb{R}^D, y_n \in \mathbb{R})\}_{n=1}^N$
- Key different from classification
 - The **loss function** is different!
 - Leading to different models and algorithms

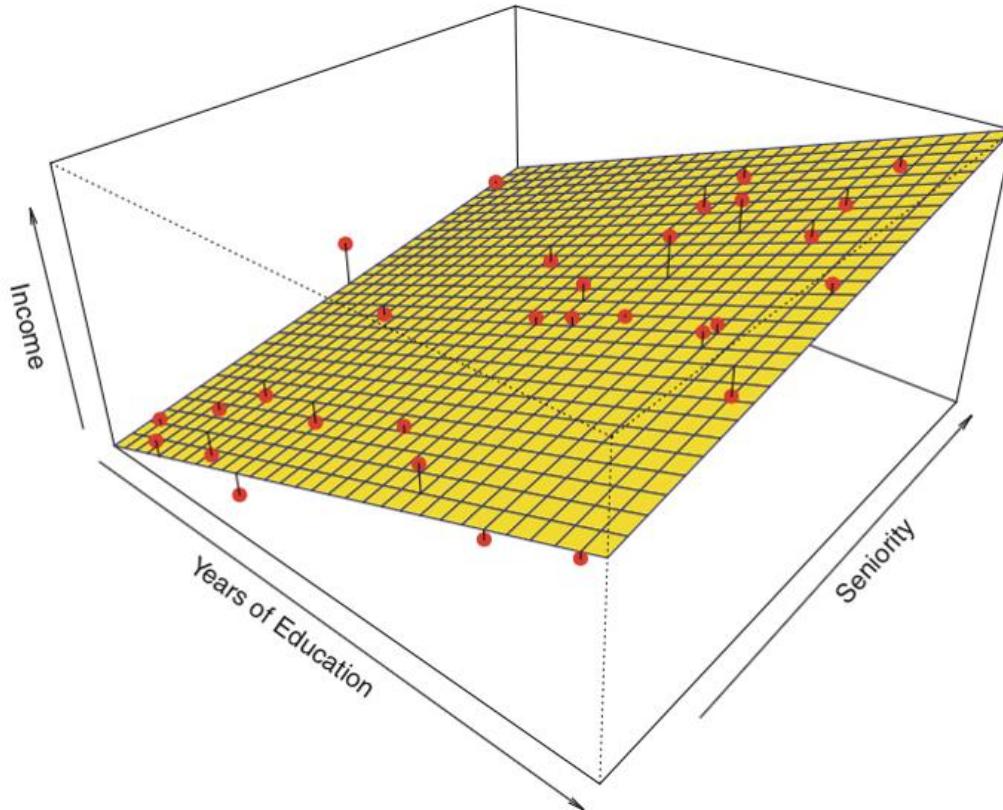
Linear regression

- Data visualization and task illustration
- Meaning:
 - Sale price = `price-per-sqft` * sqft + `fixed_expense`
- How to measure the loss?
 - Square loss: $\ell(y, \hat{y}) = (y - \hat{y})^2$
 - Absolute loss: $\ell(y, \hat{y}) = |y - \hat{y}|$

Regression (house price):
From x (area), predict y (price)



Linear regression: beyond 1-D input



Linear regression

- Training data: $D_{tr} = \{(x_n \in \mathbb{R}^D, y_n \in \mathbb{R})\}_{n=1}^N$
- Models (hypothesis class):
 - $\hat{y} = h(x) = w[1]x[1] + \dots + w[D]x[D] + b = \sum_{d=1}^D w[d]x[d] + b = \mathbf{w}^T \mathbf{x} + b$
 - $w \in \mathbb{R}^D$: weights, parameters, parameter weights
 - b (or usually denoted by $w[0]$): bias, intercept
 - Concatenation: $\tilde{w} = [w[0], w[1], \dots, w[D]]^T \in \mathbb{R}^{D+1}$; $\tilde{x} = [\mathbf{1}, x[1], \dots, x[D]]^T \in \mathbb{R}^{D+1}$
 - $\mathbf{w}^T \mathbf{x} + w[0] = \tilde{w}^T \tilde{x}$
 - Pay attention to the context: People may use notations differently!

Linear regression

- Models (hypothesis class):
 - $\hat{y} = h(\mathbf{x}) = \mathbf{w}[1]\mathbf{x}[1] + \cdots + \mathbf{w}[D]\mathbf{x}[D] + \mathbf{b} = \sum_{d=1}^D \mathbf{w}[d]\mathbf{x}[d] + \mathbf{b} = \mathbf{w}^T \mathbf{x} + \mathbf{b}$
- Assumptions:
 - Linear relationship between input features and output label
 - Linear combinations of input features
 - Parameters are coefficients of each input feature dimension $\mathbf{x}[d]$
 - Thus, only works when we expect the input-output relationship is a line/plane/hyperplane

Learning goal

- Ideally, minimizing the generalization error as much as possible
 - But we only have training data
 - So, let's try empirical risk minimization
- **Sum of Squared Error (SSE) or Residual sum of square (RSS)**
 - $\text{RSS}(\tilde{\mathbf{w}}) = \sum_i [y_i - h(\mathbf{x}_i)]^2 = \sum_i [y_i - (\mathbf{w}^T \mathbf{x}_i + b)]^2 = \sum_i [y_i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i]^2$
 - Simplicity
 - Gaussian noise assumption (later)
- Others:
 - **Sum of Absolute Error (SAE):** $\sum_i |y_i - h(\mathbf{x}_i)|$
 - **Maximum error:** $\max_i |y_i - h(\mathbf{x}_i)|$

Questions?



THE OHIO STATE UNIVERSITY

Least mean square (LMS) solution

- $\tilde{\mathbf{w}}^{\text{LMS}} = \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \text{RSS}(\tilde{\mathbf{w}}) = \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_i [y_i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i]^2$
- Let's try to expand it!

$$\begin{aligned}\text{RSS}(\tilde{\mathbf{w}}) &= \sum_i (y_i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i)^2 \\ &= \sum_i (y_i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i)(y_i - \tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}}) \\ &= \sum_i y_i^2 + \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} - 2y_i \tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} = \sum_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} - 2y_i \tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} + \text{constant}\end{aligned}$$

Least mean square (LMS) solution

- Data notation

$$X \in \mathbb{R}^{D \times N} = [x_1, \quad x_2, \quad x_3, \quad \dots, \quad x_N] = \begin{bmatrix} | & | & | & & | \\ x_1 & x_2 & x_3 & \dots & x_N \\ | & | & | & & | \end{bmatrix}$$

$$\tilde{X} \in \mathbb{R}^{(D+1) \times N} = [\tilde{x}_1, \quad \tilde{x}_2, \quad \tilde{x}_3, \quad \dots, \quad \tilde{x}_N] = \begin{bmatrix} | & | & | & & | \\ \tilde{x}_1 & \tilde{x}_2 & \tilde{x}_3 & \dots & \tilde{x}_N \\ | & | & | & & | \end{bmatrix}$$

$$y \in \mathbb{R}^N = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

Least mean square (LMS) solution

- Continued

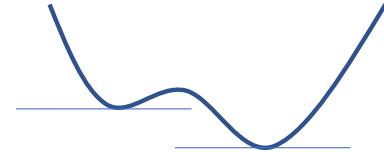
$$= \sum_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} - 2y_i \tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} + \text{constant}$$

$$= \tilde{\mathbf{w}}^T \left(\sum_i \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \right) \tilde{\mathbf{w}} - 2 \left(\sum_i y_i \tilde{\mathbf{x}}_i^T \right) \tilde{\mathbf{w}} + \text{constant}$$

$$= \tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - 2(\tilde{\mathbf{X}} \mathbf{y})^T \tilde{\mathbf{w}} + \text{constant}$$

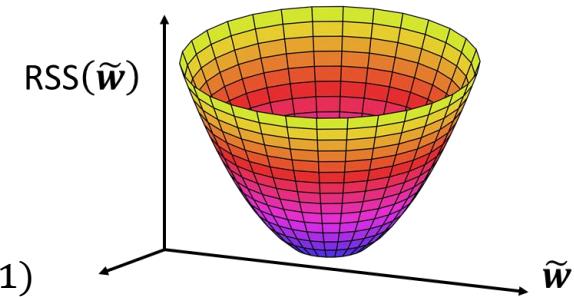
$$\boxed{\mathbf{X} \mathbf{X}^T = \begin{bmatrix} | & | & | & & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \mathbf{x}_3 & \dots & \mathbf{x}_N \\ | & | & | & & | \end{bmatrix} \begin{bmatrix} -\mathbf{x}_1^T & - \\ \vdots & \\ -\mathbf{x}_N^T & - \end{bmatrix} = \sum_{i=1}^N \begin{bmatrix} | \\ \mathbf{x}_i \\ | \end{bmatrix} \begin{bmatrix} -\mathbf{x}_i^T & - \end{bmatrix}}$$

Least mean square (LMS) solution



- Normal equations
 - Taking derivative of $\text{RSS}(\tilde{\mathbf{w}})$ w.r.t. $\tilde{\mathbf{w}}$, and setting it to **zeros**

$$\frac{\partial \text{RSS}(\tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}}} = \frac{\partial (\tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - 2(\tilde{\mathbf{X}}\mathbf{y})^T \tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}}} = 2\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T\tilde{\mathbf{w}} - 2\tilde{\mathbf{X}}\mathbf{y} = \mathbf{0} \in \mathbb{R}^{(D+1)}$$



Closed-form $\tilde{\mathbf{w}}^{\text{LMS}} = (\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1}\tilde{\mathbf{X}}\mathbf{y}$, since $2\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T\tilde{\mathbf{w}} = 2\tilde{\mathbf{X}}\mathbf{y}$ leads to $\tilde{\mathbf{w}} = (\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1}\tilde{\mathbf{X}}\mathbf{y}$

$$E(\mathbf{w}) = \mathbf{x}^T \mathbf{w}: \quad \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = \begin{bmatrix} \frac{\partial E(\mathbf{w})}{\partial w[1]} \\ \vdots \\ \frac{\partial E(\mathbf{w})}{\partial w[D]} \end{bmatrix} = \mathbf{x}$$

$$E(\mathbf{w}) = \mathbf{w}^T \mathbf{A} \mathbf{w}: \quad \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = \mathbf{A} \mathbf{w} + \mathbf{A}^T \mathbf{w} = (\mathbf{A} + \mathbf{A}^T) \mathbf{w}$$

More useful tools in the Matrix Cookbook!

Least mean square (LMS) solution: another way

$$\text{RSS}(\tilde{\mathbf{w}})$$

$$= \sum_i (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i - y_i)^2$$

$$= \sum_i (\tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} - y_i)^2$$

$$= [(\tilde{\mathbf{x}}_1^T \tilde{\mathbf{w}} - y_1), \dots, (\tilde{\mathbf{x}}_N^T \tilde{\mathbf{w}} - y_N)] \begin{bmatrix} \tilde{\mathbf{x}}_1^T \tilde{\mathbf{w}} - y_1 \\ \vdots \\ \tilde{\mathbf{x}}_N^T \tilde{\mathbf{w}} - y_N \end{bmatrix}$$

$$= \|\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - \mathbf{y}\|_2^2$$

$$= (\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - \mathbf{y})^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - \mathbf{y})$$

$$= \tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - 2(\tilde{\mathbf{X}} \mathbf{y})^T \tilde{\mathbf{w}} + \text{constant}$$

Mini summary

- Least mean square (LMS) has a closed-form solution
- Works for any model with linear combination of parameters to minimize RSS
- Bottleneck
 - $(\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1}$: computationally $O((D + 1)^3)$
 - $(\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1}$: exists only if $\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T$ is of full rank: $N \geq D + 1$. In other words, we need at least $D+1$ training data
 - Often, we need more than that to overcome noise and over-fitting

Questions?



THE OHIO STATE UNIVERSITY

Recap: Eigen-decomposition (ED)

- Given a square matrix $A \in R^{D \times D}$
- If A has D linearly independent eigenvectors $\{q_1, \dots, q_D\}$ and the corresponding eigenvalues $\{\lambda_1, \dots, \lambda_D\}$

$$A \begin{bmatrix} | & & | \\ q_1 & \dots & q_D \\ | & & | \end{bmatrix} = AQ = \begin{bmatrix} | & & | \\ \lambda_1 q_1 & \dots & \lambda_D q_D \\ | & & | \end{bmatrix} = \begin{bmatrix} | & & | \\ q_1 & \dots & q_D \\ | & & | \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \lambda_D \end{bmatrix} = Q\Lambda$$

- Eigen-decomposition: $A = Q\Lambda Q^{-1}$
 - $(AQ)Q^{-1} = (Q\Lambda)Q^{-1}$ then $A = Q\Lambda Q^{-1}$
- Not all the square matrices have eigen-decomposition

Recap: Singular value decomposition (SVD)

- Singular value decomposition (SVD)

Singular value decomposition (SVD)

$$X = U\Gamma V^T = \begin{bmatrix} | & | & | & & | \\ u_1 & u_2 & u_3 & \dots & u_D \\ | & | & | & & | \end{bmatrix} \begin{bmatrix} \gamma_1 & 0 & 0 & 0 \\ 0 & \ddots & 0 & 0 \\ 0 & 0 & \gamma_D & 0 \end{bmatrix} \begin{bmatrix} -v_1^T & - \\ \vdots & \\ -v_N^T & - \end{bmatrix}$$

Non-negative

$$\mathbf{U}^T \mathbf{U} = \begin{bmatrix} -\mathbf{u}_1^T & - \\ \vdots & \\ -\mathbf{u}_N^T & - \end{bmatrix} \begin{bmatrix} | & | & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \mathbf{u}_3 & \dots & \mathbf{u}_D \\ | & | & | & & | \end{bmatrix} = \mathbf{I} = \mathbf{U} \mathbf{U}^T \quad \text{Orthonormal, unitary}$$

$$V^T V = \begin{bmatrix} -\boldsymbol{v}_1^T & - \\ \vdots & \\ -\boldsymbol{v}_N^T & - \end{bmatrix} \begin{bmatrix} | & | & | & & | \\ \boldsymbol{v}_1 & \boldsymbol{v}_2 & \boldsymbol{v}_3 & \dots & \boldsymbol{v}_D \\ | & | & | & & | \end{bmatrix} = I = V V^T \quad \text{Orthonormal, unitary}$$

Recap: SVD + ED

- Combination (suppose $X \in R^{D \times N}$ and $X = U\Gamma V^T$)

$$\begin{aligned} XX^T &= U\Gamma V^T (U\Gamma V^T)^T = U\Gamma V^T (V\Gamma^T U^T) = U\Gamma V^T V\Gamma^T U^T = U\Gamma\Gamma^T U^T \\ &= U(\Gamma\Gamma^T)U^T = U(\Gamma\Gamma^T)U^{-1} \end{aligned}$$

- XX^T has an eigen-decomposition $U(\Gamma\Gamma^T)U^{-1}$

Definition of SVD

- Inversion

$$(XX^T)^{-1} = (U(\Gamma\Gamma^T)U^{-1})^{-1} = U(\Gamma\Gamma^T)^{-1}U^{-1}$$

- XX^T is invertible iff $\Gamma\Gamma^T$ is invertible

Why not invertible?

- Look deeper!

$$\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T = \mathbf{U} \begin{bmatrix} \gamma_1^2 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \gamma_D^2 \end{bmatrix} \mathbf{U}^T$$

$$(\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1} = \mathbf{U} \begin{bmatrix} \frac{1}{\gamma_1^2} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \frac{1}{\gamma_D^2} \end{bmatrix} \mathbf{U}^T$$

- Not invertible if there are zeros along the diagonal

How to make it invertible?

- Adding positive numbers to the diagonal!

$$\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T + \lambda \mathbf{I} = \mathbf{U} \begin{bmatrix} \gamma_1^2 + \lambda & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \gamma_D^2 + \lambda \end{bmatrix} \mathbf{U}^T \quad (\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T + \lambda \mathbf{I})^{-1} = \mathbf{U} \begin{bmatrix} \frac{1}{\gamma_1^2 + \lambda} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \frac{1}{\gamma_D^2 + \lambda} \end{bmatrix} \mathbf{U}^T$$

- Closed-form solution = $(\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T + \lambda \mathbf{I})^{-1}\tilde{\mathbf{X}}\mathbf{y}$
 - Named “ridge regression”
 - Not the same $\tilde{\mathbf{w}}^{\text{LMS}}$ when $\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T$ is invertible
- We will justify why doing this makes sense later!

Practice!

- Show $(\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T + \lambda \mathbf{I})^{-1}\tilde{\mathbf{X}}\mathbf{y}$ is the solution of $\underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_i [y_i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i]^2 + \lambda \|\tilde{\mathbf{w}}\|_2^2$
 - That is, we introduce the prior knowledge that the solution should be close to $\mathbf{0}$
- How to select λ ?

Summary

- Linear regression (linear curve fitting)
 - Model input-output relationships via linear functions
- To learn regression model parameters
 - Given training data examples
 - Choose an error/loss function
- Linear least mean square (LMS)
 - Minimize sum of residual sum of square (RSS)
 - Calculus and Linear Algebra
 - Ridge regression for noninvertible cases

CSE 5523: Regression



THE OHIO STATE UNIVERSITY

Hw-1

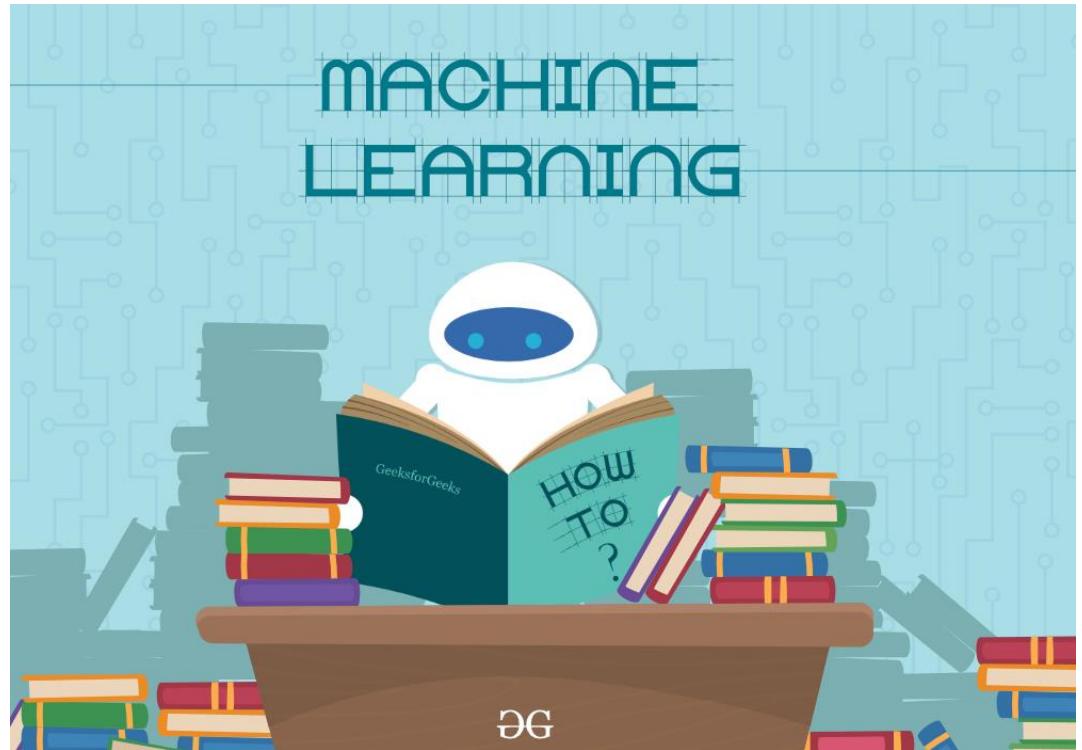
- Will be released tonight (midnight)
 - You will receive Carmen's announcement
- Will be **due on 9/17** (midnight)
- Questions in both the programming set and the problem set
- Please use Piazza or come to the office hours for discussion
- Please start working on it ASAP

Today

Linear regression

- Review
- Linear regression (continued)
- Ridge regression
- A little more math

Nonlinear regression



Linear regression

- Training data: $D_{tr} = \{(x_n \in \mathbb{R}^D, y_n \in \mathbb{R})\}_{n=1}^N$
- Models (hypothesis class):
 - $\hat{y} = h(x) = w^T x + w[0] = \tilde{w}^T \tilde{x}$
 - $\tilde{w} = [w[0], w[1], \dots, w[D]]^T \in \mathbb{R}^{D+1}; \tilde{x} = [1, x[1], \dots, x[D]]^T \in \mathbb{R}^{D+1}$
 - Assumption: Linear relationship between input features and output label

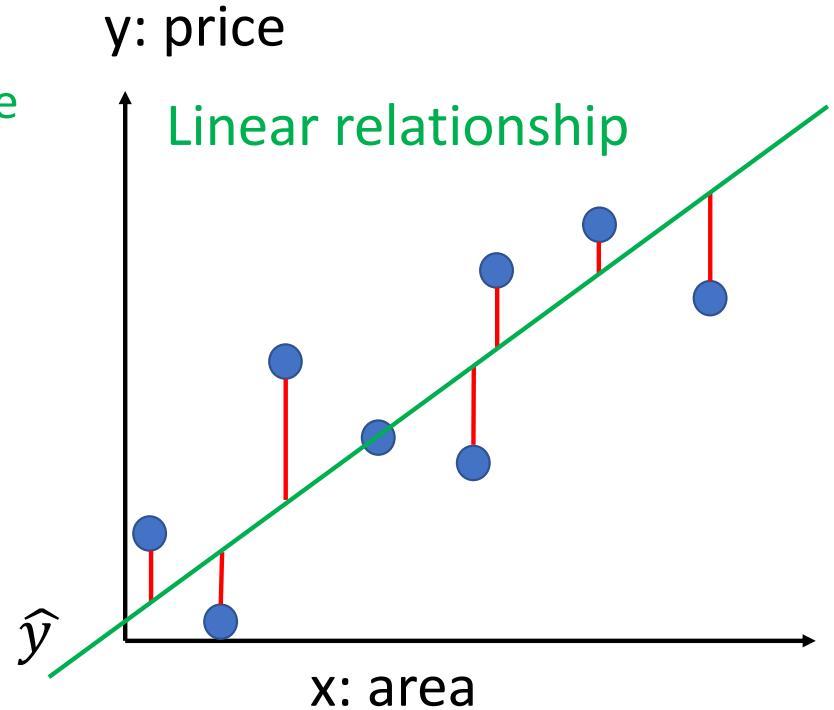
$$X \in \mathbb{R}^{D \times N} = [x_1, x_2, x_3, \dots, x_N] = \begin{bmatrix} | & | & | & \dots & | \\ x_1 & x_2 & x_3 & \dots & x_N \\ | & | & | & & | \end{bmatrix} \quad y \in \mathbb{R}^N = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

$$\tilde{X} \in \mathbb{R}^{(D+1) \times N} = [\tilde{x}_1, \tilde{x}_2, \tilde{x}_3, \dots, \tilde{x}_N] = \begin{bmatrix} | & | & | & \dots & | \\ \tilde{x}_1 & \tilde{x}_2 & \tilde{x}_3 & \dots & \tilde{x}_N \\ | & | & | & & | \end{bmatrix}$$

Linear regression

- Data visualization and task illustration
- Meaning:
 - Sale price = **price-per-sqft** * sqft + **fixed_expense**
- Learning goal?
 - Estimate parameters by minimizing the loss
 - **Residual sum of square (RSS):**
 - $\text{RSS}(\tilde{w}) = \sum_i [y_i - h(x_i)]^2$

Regression (house price):
From x (area), predict y (price)

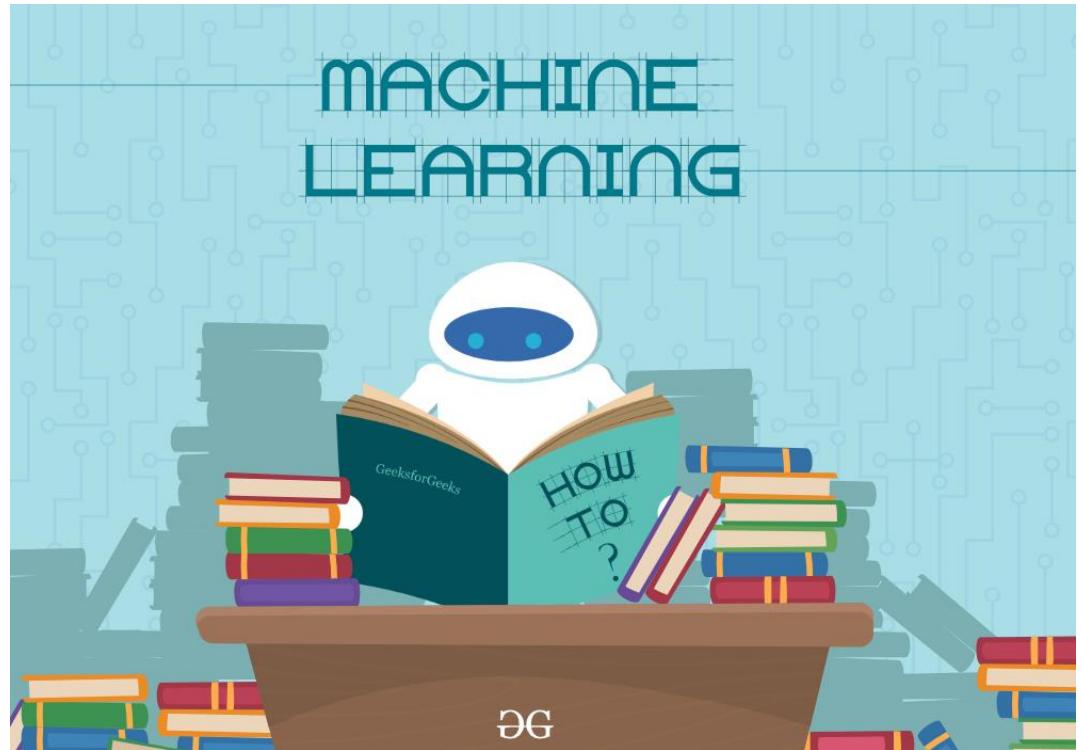


Today

Linear regression

- Review
- Linear regression (continued)
- Ridge regression
- A little more math

Nonlinear regression



Least mean square (LMS) solution

- $\tilde{\mathbf{w}}^{\text{LMS}} = \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \text{RSS}(\tilde{\mathbf{w}}) = \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_i [y_i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i]^2$
- Let's try to expand it!

$$\begin{aligned}\text{RSS}(\tilde{\mathbf{w}}) &= \sum_i (y_i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i)^2 \\ &= \sum_i (y_i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i)(y_i - \tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}}) \\ &= \sum_i y_i^2 + \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} - 2y_i \tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} = \sum_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} - 2y_i \tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} + \text{constant}\end{aligned}$$

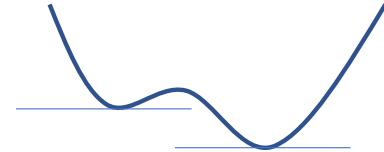
Least mean square (LMS) solution

- Continued

$$\begin{aligned}\text{RSS}(\tilde{\mathbf{w}}) &= \sum_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} - 2y_i \tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} + \text{constant} \\ &= \tilde{\mathbf{w}}^T \left(\sum_i \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^T \right) \tilde{\mathbf{w}} - 2 \left(\sum_i y_i \tilde{\mathbf{x}}_i^T \right) \tilde{\mathbf{w}} + \text{constant} \\ &= \tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - 2(\tilde{\mathbf{X}} \mathbf{y})^T \tilde{\mathbf{w}} + \text{constant}\end{aligned}$$

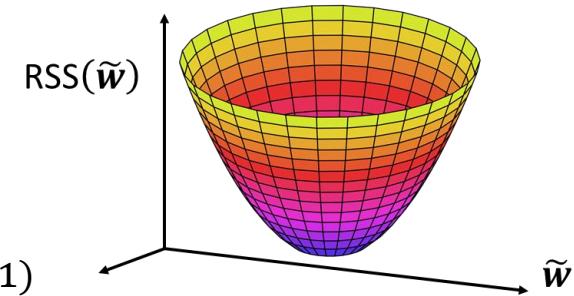
$$\boxed{\mathbf{X} \mathbf{X}^T = \begin{bmatrix} | & | & | & & | \\ \mathbf{x}_1 & \mathbf{x}_2 & \mathbf{x}_3 & \dots & \mathbf{x}_N \\ | & | & | & & | \end{bmatrix} \begin{bmatrix} -\mathbf{x}_1^T & - \\ \vdots & \\ -\mathbf{x}_N^T & - \end{bmatrix} = \sum_{i=1}^N \begin{bmatrix} | \\ \mathbf{x}_i \\ | \end{bmatrix} \begin{bmatrix} -\mathbf{x}_i^T & - \end{bmatrix}}$$

Least mean square (LMS) solution



- Normal equations
 - Taking derivative of $\text{RSS}(\tilde{\mathbf{w}})$ w.r.t. $\tilde{\mathbf{w}}$, and setting it to zeros

$$\frac{\partial \text{RSS}(\tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}}} = \frac{\partial (\tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - 2(\tilde{\mathbf{X}}\mathbf{y})^T \tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}}} = 2\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T\tilde{\mathbf{w}} - 2\tilde{\mathbf{X}}\mathbf{y} = \mathbf{0} \in \mathbb{R}^{(D+1)}$$



Closed-form $\tilde{\mathbf{w}}^{\text{LMS}} = (\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1}\tilde{\mathbf{X}}\mathbf{y}$, since $2\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T\tilde{\mathbf{w}} = 2\tilde{\mathbf{X}}\mathbf{y}$ leads to $\tilde{\mathbf{w}} = (\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1}\tilde{\mathbf{X}}\mathbf{y}$

$$E(\mathbf{w}) = \mathbf{x}^T \mathbf{w}: \quad \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = \begin{bmatrix} \frac{\partial E(\mathbf{w})}{\partial w[1]} \\ \vdots \\ \frac{\partial E(\mathbf{w})}{\partial w[D]} \end{bmatrix} = \mathbf{x}$$

$$E(\mathbf{w}) = \mathbf{w}^T \mathbf{A} \mathbf{w}: \quad \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = \mathbf{A} \mathbf{w} + \mathbf{A}^T \mathbf{w} = (\mathbf{A} + \mathbf{A}^T) \mathbf{w}$$

More useful tools in the Matrix Cookbook!

Least mean square (LMS) solution: another way

$$\text{RSS}(\tilde{\mathbf{w}})$$

$$= \sum_i (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i - y_i)^2$$

$$= \sum_i (\tilde{\mathbf{x}}_i^T \tilde{\mathbf{w}} - y_i)^2$$

$$= [(\tilde{\mathbf{x}}_1^T \tilde{\mathbf{w}} - y_1), \dots, (\tilde{\mathbf{x}}_N^T \tilde{\mathbf{w}} - y_N)] \begin{bmatrix} \tilde{\mathbf{x}}_1^T \tilde{\mathbf{w}} - y_1 \\ \vdots \\ \tilde{\mathbf{x}}_N^T \tilde{\mathbf{w}} - y_N \end{bmatrix}$$

$$= \|\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - \mathbf{y}\|_2^2$$

$$= (\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - \mathbf{y})^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - \mathbf{y})$$

$$= \tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - 2(\tilde{\mathbf{X}} \mathbf{y})^T \tilde{\mathbf{w}} + \text{constant}$$

Mini summary

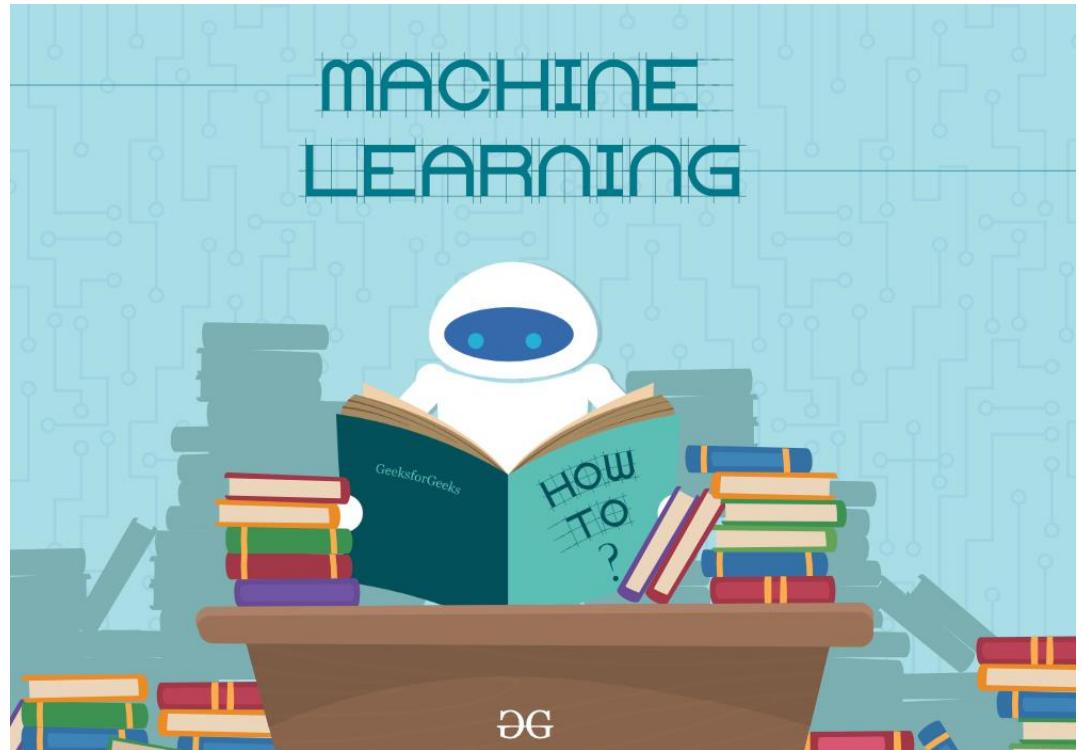
- Least mean square (LMS) has a closed-form solution
- Works for any model with linear combination of parameters to minimize RSS
- Bottleneck
 - $(\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1}$: computationally $O((D + 1)^3)$
 - $(\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1}$: exists only if $\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T$ is of full rank: $N \geq D + 1$. In other words, we need at least $D+1$ training data
 - Often, we need more than that to overcome noise and over-fitting

Today

Linear regression

- Review
- Linear regression (continued)
- Ridge regression
- A little more math

Nonlinear regression



Recap: Eigen-decomposition (ED)

- Given a square matrix $A \in R^{D \times D}$
- If A has D linearly independent eigenvectors $\{q_1, \dots, q_D\}$ and the corresponding eigenvalues $\{\lambda_1, \dots, \lambda_D\}$

$$A \begin{bmatrix} | & & | \\ q_1 & \dots & q_D \\ | & & | \end{bmatrix} = AQ = \begin{bmatrix} | & & | \\ \lambda_1 q_1 & \dots & \lambda_D q_D \\ | & & | \end{bmatrix} = \begin{bmatrix} | & & | \\ q_1 & \dots & q_D \\ | & & | \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \lambda_D \end{bmatrix} = Q\Lambda$$

- Eigen-decomposition: $A = Q\Lambda Q^{-1}$
 - $(AQ)Q^{-1} = (Q\Lambda)Q^{-1}$ then $A = Q\Lambda Q^{-1}$
- Not all the square matrices have eigen-decomposition

Recap: Singular value decomposition (SVD)

- Singular value decomposition (SVD)

Singular value decomposition (SVD)

$$X = U\Gamma V^T = \begin{bmatrix} | & | & | & & | \\ u_1 & u_2 & u_3 & \dots & u_D \\ | & | & | & & | \end{bmatrix} \begin{bmatrix} \gamma_1 & 0 & 0 & 0 \\ 0 & \ddots & 0 & 0 \\ 0 & 0 & \gamma_D & 0 \end{bmatrix} \begin{bmatrix} -v_1^T & - \\ \vdots & \\ -v_N^T & - \end{bmatrix}$$

Non-negative

$$\mathbf{U}^T \mathbf{U} = \begin{bmatrix} -\mathbf{u}_1^T & - \\ \vdots & \\ -\mathbf{u}_N^T & - \end{bmatrix} \begin{bmatrix} | & | & | & & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \mathbf{u}_3 & \dots & \mathbf{u}_D \\ | & | & | & & | \end{bmatrix} = \mathbf{I} = \mathbf{U} \mathbf{U}^T \quad \text{Orthonormal, unitary}$$

$$V^T V = \begin{bmatrix} -\boldsymbol{v}_1^T & - \\ \vdots & \\ -\boldsymbol{v}_N^T & - \end{bmatrix} \begin{bmatrix} | & | & | & & | \\ \boldsymbol{v}_1 & \boldsymbol{v}_2 & \boldsymbol{v}_3 & \dots & \boldsymbol{v}_D \\ | & | & | & & | \end{bmatrix} = I = VV^T \quad \text{Orthonormal, unitary}$$

Recap: SVD + ED

- Combination (suppose $X \in R^{D \times N}$ and $X = U\Gamma V^T$)

$$\begin{aligned} XX^T &= U\Gamma V^T (U\Gamma V^T)^T = U\Gamma V^T (V\Gamma^T U^T) = U\Gamma V^T V\Gamma^T U^T = U\Gamma\Gamma^T U^T \\ &= U(\Gamma\Gamma^T)U^T = U(\Gamma\Gamma^T)U^{-1} \end{aligned}$$

- XX^T has an eigen-decomposition $U(\Gamma\Gamma^T)U^{-1}$

Definition of SVD

- Inversion

$$(XX^T)^{-1} = (U(\Gamma\Gamma^T)U^{-1})^{-1} = U(\Gamma\Gamma^T)^{-1}U^{-1}$$

- XX^T is invertible iff $\Gamma\Gamma^T$ is invertible

Why not invertible?

- Look deeper!

$$\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T = \mathbf{U} \begin{bmatrix} \gamma_1^2 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \gamma_D^2 \end{bmatrix} \mathbf{U}^T$$

$$(\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1} = \mathbf{U} \begin{bmatrix} \frac{1}{\gamma_1^2} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \frac{1}{\gamma_D^2} \end{bmatrix} \mathbf{U}^T$$

- Not invertible if there are zeros along the diagonal

How to make it invertible?

- Adding positive numbers to the diagonal!

$$\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T + \lambda \mathbf{I} = \mathbf{U} \begin{bmatrix} \gamma_1^2 + \lambda & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \gamma_D^2 + \lambda \end{bmatrix} \mathbf{U}^T \quad (\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T + \lambda \mathbf{I})^{-1} = \mathbf{U} \begin{bmatrix} \frac{1}{\gamma_1^2 + \lambda} & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \frac{1}{\gamma_D^2 + \lambda} \end{bmatrix} \mathbf{U}^T$$

- Closed-form solution = $(\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T + \lambda \mathbf{I})^{-1}\tilde{\mathbf{X}}y$
 - Named “ridge regression”
 - Not the same $\tilde{\mathbf{w}}^{\text{LMS}}$ when $\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T$ is invertible
- We will justify why doing this makes sense later!

Practice!

- Show $(\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T + \lambda \mathbf{I})^{-1}\tilde{\mathbf{X}}\mathbf{y}$ is the solution of $\underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_i [y_i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i]^2 + \lambda \|\tilde{\mathbf{w}}\|_2^2$
 - That is, we introduce the prior knowledge that the solution should be close to $\mathbf{0}$
- How to select λ ?

Summary

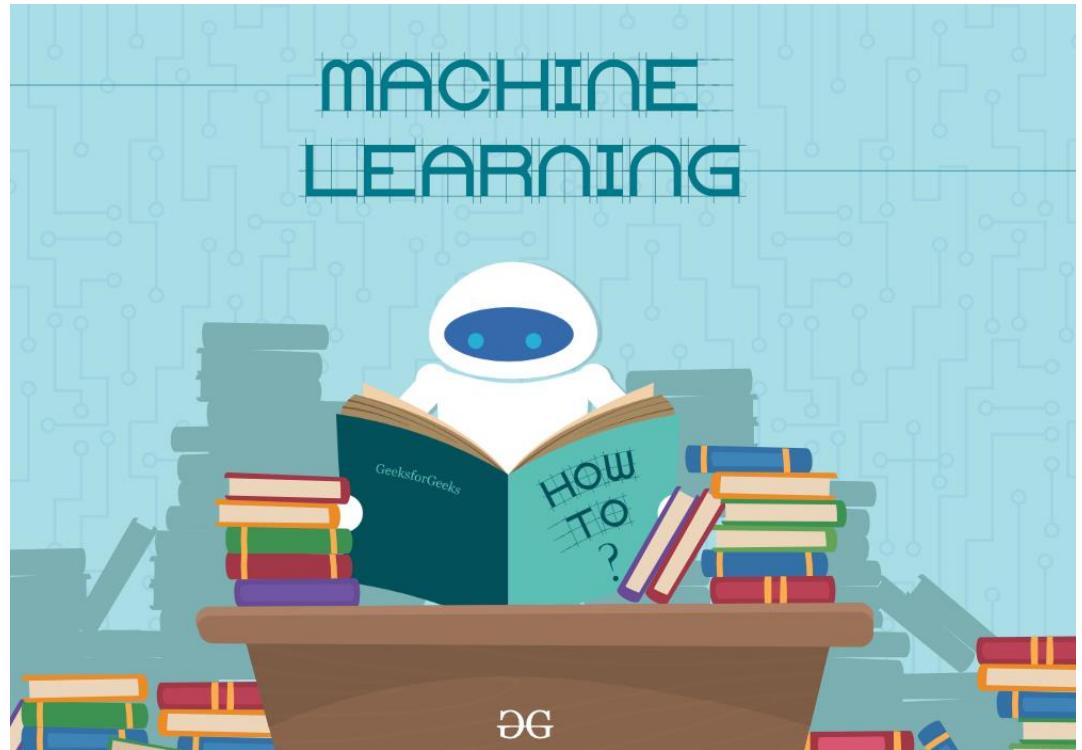
- Linear regression (linear curve fitting)
 - Model input-output relationships via linear functions
- To learn regression model parameters
 - Given training data examples
 - Choose an error/loss function
- Linear least mean square (LMS)
 - Minimize sum of residual sum of square (RSS)
 - Calculus and Linear Algebra
 - Ridge regression for noninvertible cases

Today

Linear regression

- Review
- Linear regression (continued)
- Ridge regression
- A little more math

Nonlinear regression



A little more math

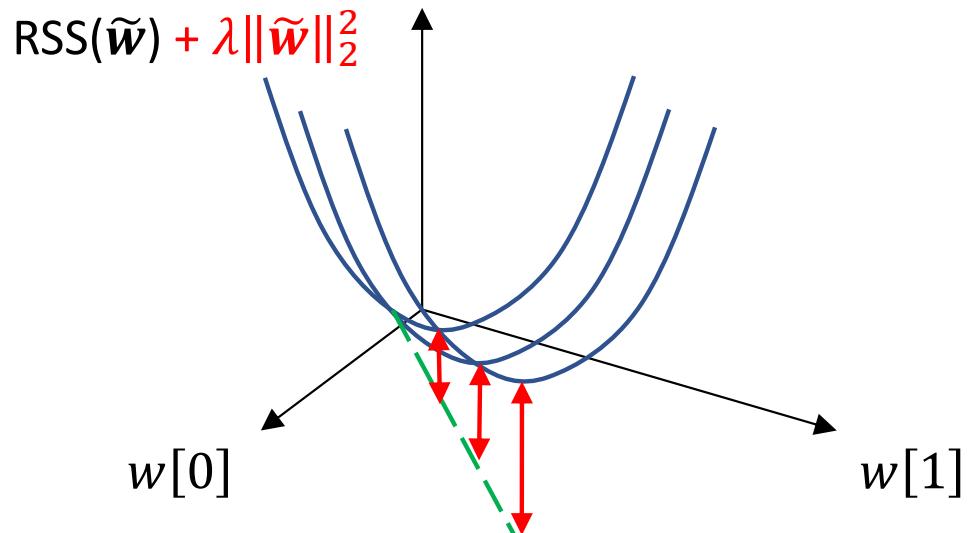
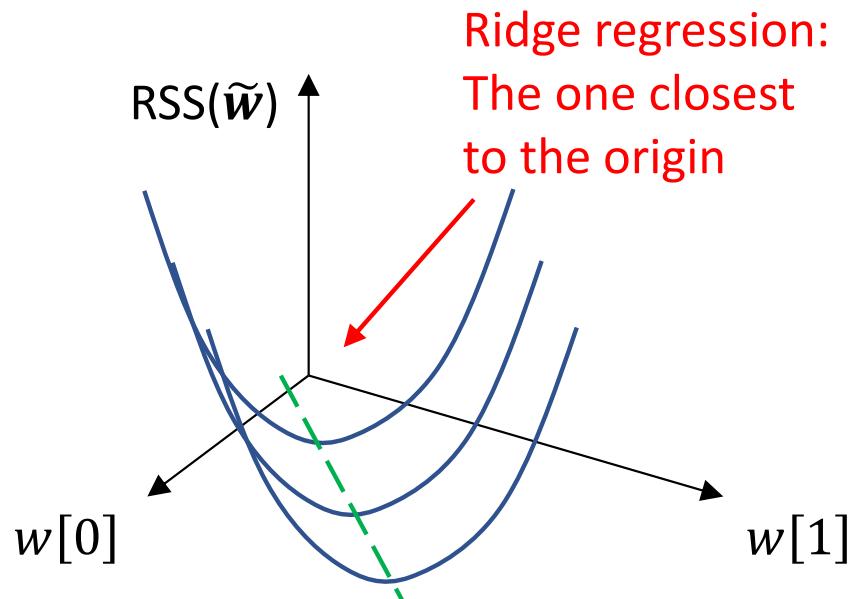
- Goal:
 - We want to minimize $\text{RSS}(\tilde{\mathbf{w}}) = \sum_i [y_i - h(\mathbf{x}_i)]^2$
 - Normal equations tell us $\frac{\partial \text{RSS}(\tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}}} = 2\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T\tilde{\mathbf{w}} - 2\tilde{\mathbf{X}}\mathbf{y} = \mathbf{0} \in \mathbb{R}^{(D+1)}$
 - Closed-form solution: $\tilde{\mathbf{w}}^{\text{LMS}} = (\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1}\tilde{\mathbf{X}}\mathbf{y}$
- What does it mean by $(\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)^{-1}$ not invertible?
 - NO solutions for $\frac{\partial \text{RSS}(\tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}}} = \mathbf{0}$?
 - TOO MANY solutions for $\frac{\partial \text{RSS}(\tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}}} = \mathbf{0}$?

A little more math

- Answer:
 - D+1 unknown (degree of freedom), but < D+1 linear independent training data in “ \tilde{X} ”
 - TOO MANY choices of unknown to fit the training data
- Examples:
 - Learn $\hat{y} = h(\mathbf{x}) = w[1]\mathbf{x} + w[0]$ from $(\mathbf{x}_1 = -1, y_1 = 1)$
 - Any pair of $(w[1], w[0])$ such that $w[1] - w[0] = -1$
 - Learn $\hat{y} = h(\mathbf{x}) = w[1]\mathbf{x} + w[0]$ from $(\mathbf{x}_1 = -1, y_1 = 1)$ and $(\mathbf{x}_2 = -1, y_2 = 1)$
 - Any pair of $(w[1], w[0])$ such that $w[1] - w[0] = -1$
 - Learn $\hat{y} = h(\mathbf{x}) = w[1]\mathbf{x} + w[0]$ from $(\mathbf{x}_1 = -1, y_1 = 1)$ and $(\mathbf{x}_2 = 1, y_2 = 1)$
 - $w[1] - w[0] = -1$ AND $w[1] + w[0] = 1$, we get $w[1] = 0, w[0] = 1$

What does ridge regression do?

- Within many solutions, identify one!
- Geometric illustration: $w[1] - w[0] = -1$



Questions?



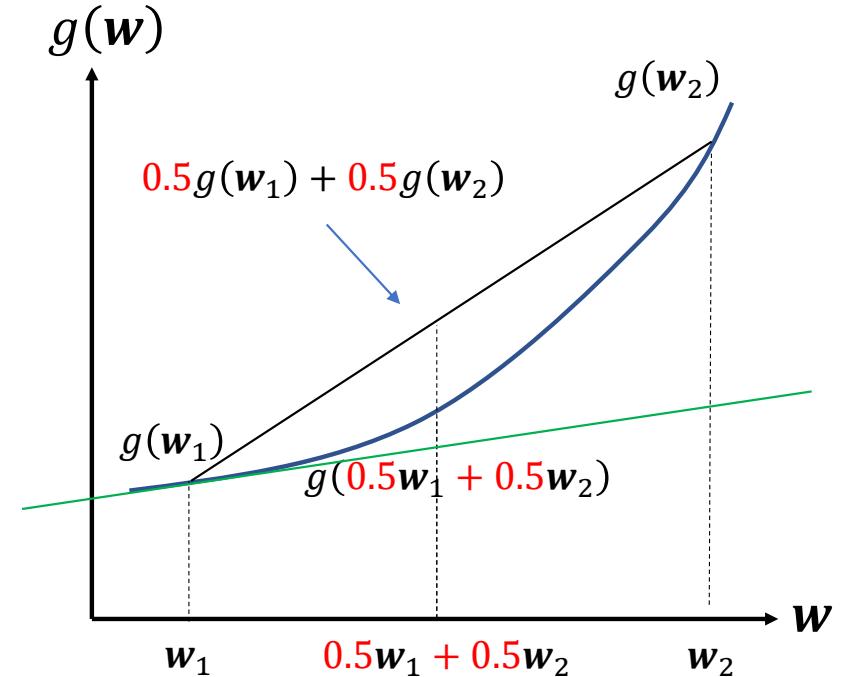
THE OHIO STATE UNIVERSITY

A little more math: convex functions

- A function $g(\mathbf{w})$ is convex iff $g(t\mathbf{w}_1 + (1 - t)\mathbf{w}_2) \leq tg(\mathbf{w}_1) + (1 - t)g(\mathbf{w}_2)$
 - $\forall \mathbf{w}_1, \mathbf{w}_2; \forall t \in [0, 1]$

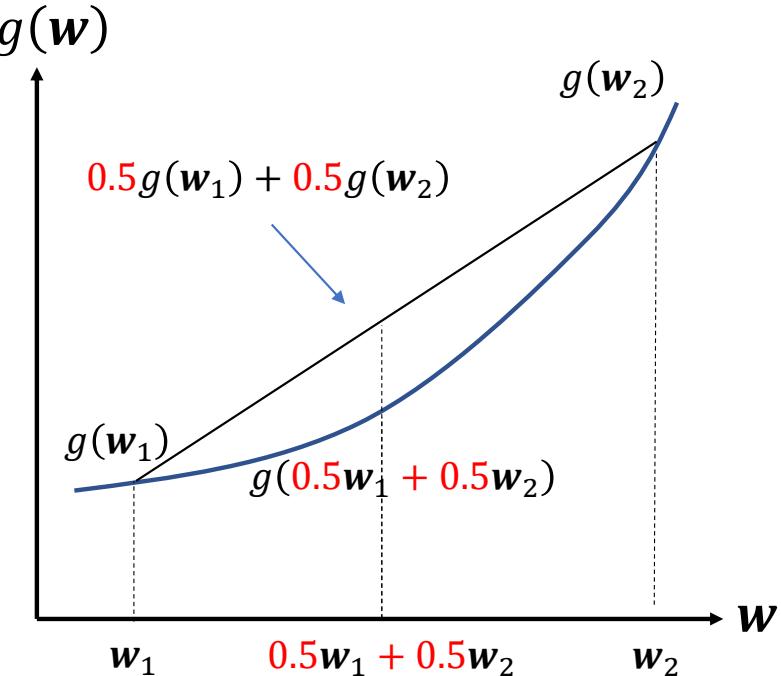
- If $g(\mathbf{w})$ is differentiable, then
 - $g(\mathbf{w}_2) \geq g(\mathbf{w}_1) + \langle \nabla g(\mathbf{w}_1), \mathbf{w}_2 - \mathbf{w}_1 \rangle$
 - $\forall \mathbf{w}_1, \mathbf{w}_2$

- Gradient $\nabla g(\mathbf{w}) = \begin{bmatrix} \frac{\partial g(\mathbf{w})}{\partial w[1]} \\ \vdots \\ \frac{\partial g(\mathbf{w})}{\partial w[D]} \end{bmatrix}$



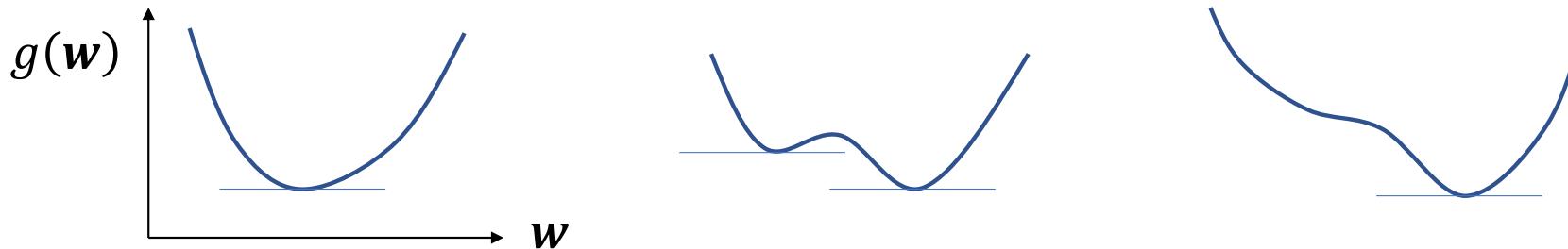
A little more math: convex functions

- A function $g(\mathbf{w})$ is convex iff $g(t\mathbf{w}_1 + (1 - t)\mathbf{w}_2) \leq tg(\mathbf{w}_1) + (1 - t)g(\mathbf{w}_2)$
 - $\forall \mathbf{w}_1, \mathbf{w}_2; \forall t \in [0, 1]$
- If $g(\mathbf{w})$ is twice differentiable, then
 - $\mathbf{H}g(\mathbf{w}) = \begin{bmatrix} \frac{\partial^2 g(\mathbf{w})}{\partial w[1]\partial w[1]} & \dots & \frac{\partial^2 g(\mathbf{w})}{\partial w[1]\partial w[D]} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 g(\mathbf{w})}{\partial w[D]\partial w[1]} & \dots & \frac{\partial^2 g(\mathbf{w})}{\partial w[D]\partial w[D]} \end{bmatrix} \geq 0, \forall \mathbf{w}$
 - Hessian matrix is positive semi-definite (PSD), $\forall \mathbf{w}$
 - Physical meaning?



A little more math: convex functions

- Properties: Any local optimum is a global optimum
 - Gradients (partial derivatives) = 0 means local optimums



- Is $\text{RSS}(\tilde{\mathbf{w}})$ convex?
 - $\text{RSS}(\tilde{\mathbf{w}}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - 2(\tilde{\mathbf{X}}\mathbf{y})^T \tilde{\mathbf{w}} + \text{constant}$
 - $\nabla \text{RSS}(\tilde{\mathbf{w}}) = 2\tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - 2\tilde{\mathbf{X}}\mathbf{y}$
 - $\mathbf{H} \text{RSS}(\tilde{\mathbf{w}}) = 2\tilde{\mathbf{X}} \tilde{\mathbf{X}}^T$
 - PSD because $\mathbf{v}^T (2\tilde{\mathbf{X}} \tilde{\mathbf{X}}^T) \mathbf{v} = 2(\tilde{\mathbf{X}}^T \mathbf{v})^T \tilde{\mathbf{X}}^T \mathbf{v} = 2\|\tilde{\mathbf{X}}^T \mathbf{v}\|_2^2 \geq 0, \forall \mathbf{v} \in \mathbb{R}^{D+1} \neq \mathbf{0}$

Caution

- When the function to be minimized is convex:
 - Any local optimum is a global optimum
 - However, it does not imply there exists closed-form solutions
 - Indeed, most of the (convex) optimization problems do not have closed-form solutions

Questions?



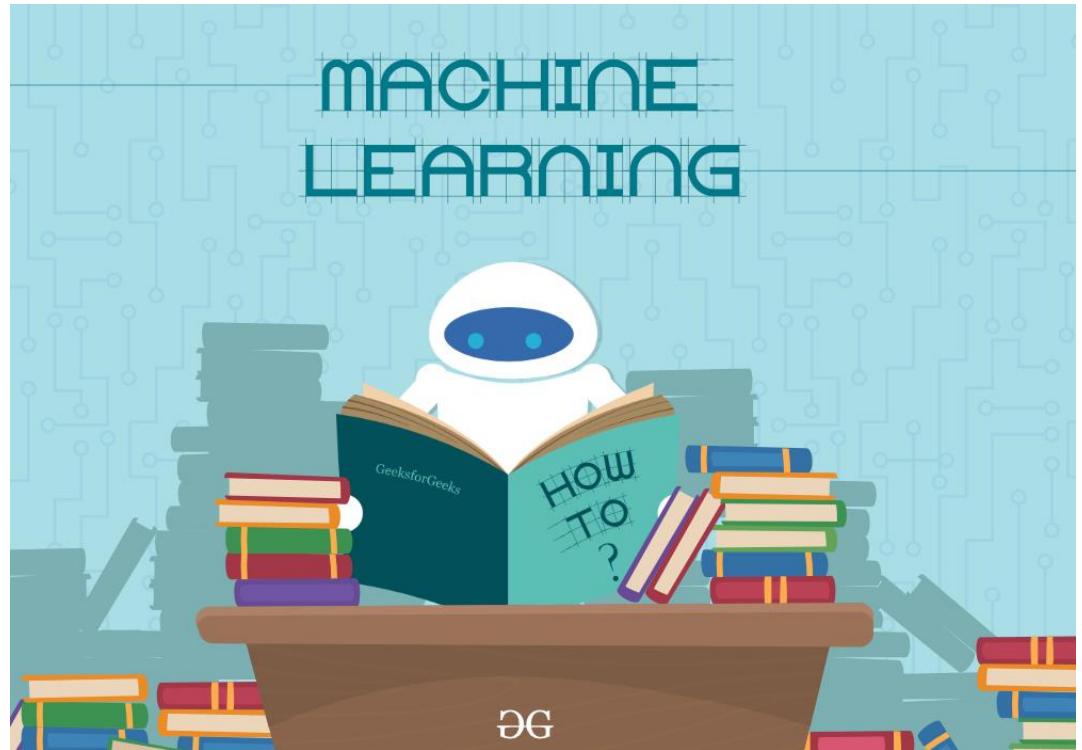
THE OHIO STATE UNIVERSITY

Today

Linear regression

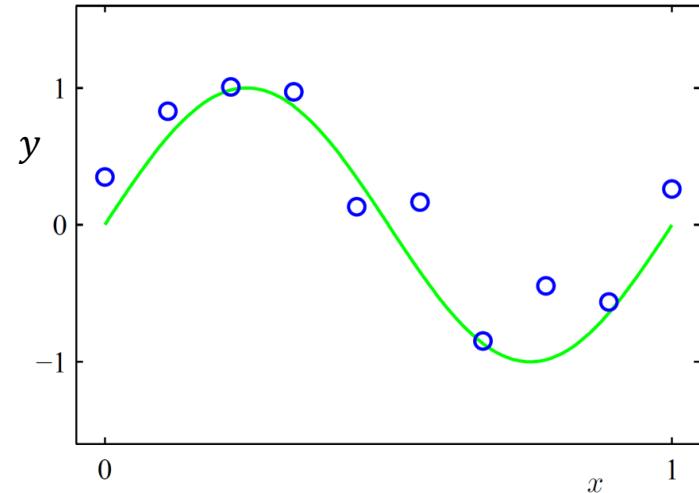
- Review
- Linear regression (continued)
- Ridge regression
- A little more math

Nonlinear regression



Nonlinear regression

- Definition: fit a **nonlinear curve** to data
 - $h(\mathbf{x})$ is linear iff $h(a\mathbf{x}_1 + b\mathbf{x}_2) = a \times h(\mathbf{x}_1) + b \times h(\mathbf{x}_2)$
 - $\forall \mathbf{x}_1, \mathbf{x}_2; \forall a, b \in \mathbb{R}$
 - Otherwise, it is nonlinear
 - Example: $h(\mathbf{x}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$ is linear
- Questions:
 - How to represent a nonlinear function?
 - What are the parameters?
 - How to estimate the parameters?



Nonlinear regression: representation

- Method 1: directly use nonlinear functions

- $h(x; w) = e^{wx}$
- $h(x; w) = \sin(wx)$
- ...

Explicitly indicate the parameters;
alternatively, $h_w(x)$

- Method 2: use nonlinear **feature transforms** or **mappings**

- $\phi(x): x \in \mathbb{R}^D \mapsto z \in \mathbb{R}^M$
- M is the dimensionality of the transformed feature z (or $\phi(x)$); can be larger, smaller than, or the same as D
- Notation $z[m] = \phi_m(x)$
- ϕ_m is called a nonlinear basis function
- $\phi(x)$ can be used for other ML algorithms (e.g., nearest neighbor) as well

Nonlinear regression: representation

- Examples of $\phi_m(x)$:
 - $z[m] = (x[m])^2$
 - $z[m] = \sin(x[m])$
 - $z[m] = \log(x[m])$
 - $z[m] = e^{x[m]}$
 - When $x \in \mathbb{R}$
 - $z[m] = (x)^{m-1}$
 - $z \in \mathbb{R}^M$ can already encode the dummy variable 1

Nonlinear regression: representation

- Previously we use $h(\mathbf{x}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$
- Now we use $h(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) = \sum_m w[m] \phi_m(\mathbf{x})$, where $\mathbf{w} \in \mathbb{R}^M$
 - In nonlinear regression, the number of parameters is usually different from D or D+1
- $\text{RSS}(\mathbf{w}) = \sum_i [y_i - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i)]^2$
- $\mathbf{w}^{\text{LMS}} = \underset{\mathbf{w}}{\operatorname{argmin}} \text{RSS}(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_i [y_i - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i)]^2$

Nonlinear regression: parameter estimation

- $\mathbf{w}^{\text{LMS}} = \underset{\mathbf{w}}{\operatorname{argmin}} \text{RSS}(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_i [y_i - \mathbf{w}^T \boldsymbol{\phi}(x_i)]^2$

$$\Phi \in \mathbb{R}^{M \times N} = [\boldsymbol{\phi}(x_1) \quad \dots \quad \boldsymbol{\phi}(x_N)] = \begin{bmatrix} | & & | \\ \boldsymbol{\phi}(x_1) & \dots & \boldsymbol{\phi}(x_N) \\ | & & | \end{bmatrix}$$

- Closed form solution \mathbf{w}^{LMS} : $(\tilde{\Phi} \tilde{\Phi}^T)^{-1} \tilde{\Phi} \mathbf{y}$

- Question:

- How many data are needed?
 - Is $\text{RSS}(\mathbf{w})$ still convex w.r.t. \mathbf{w} ?

Questions?



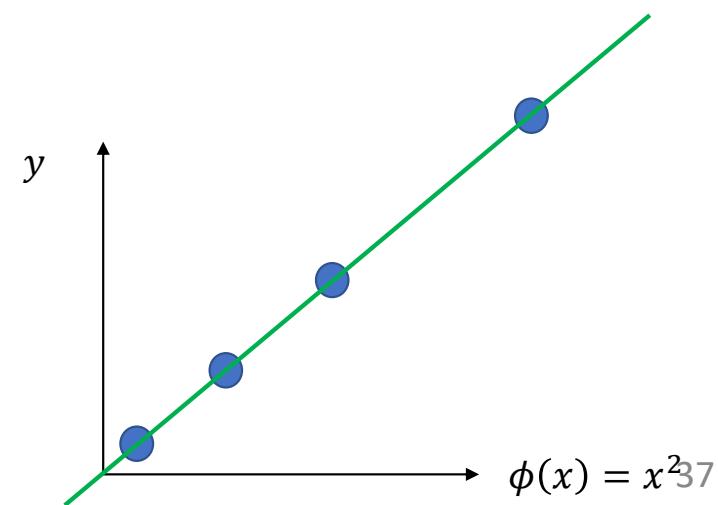
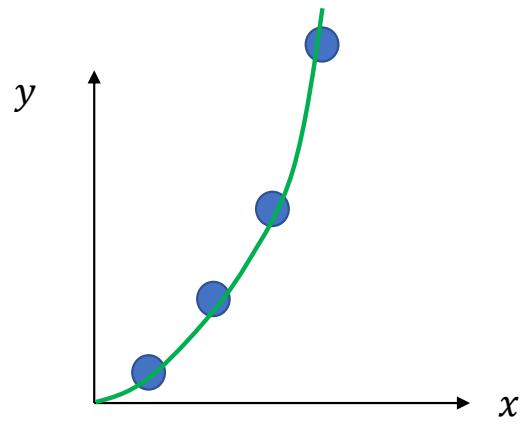
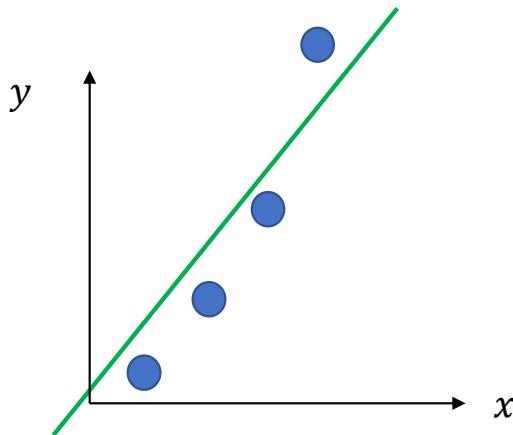
THE OHIO STATE UNIVERSITY

Nonlinear regression: example

- Polynomial regression for 1-D x , i.e., $h(x) = \sum_{m=0}^M w[m]x^m = \mathbf{w}^T \boldsymbol{\phi}(x)$

$$\boldsymbol{\phi}(x) = \begin{bmatrix} 1 \\ x \\ \vdots \\ x^M \end{bmatrix}$$

- Special case: $D_{tr} = \{(x_n \in \mathbb{R}, y_n = x_n^2)\}_{n=1}^N$

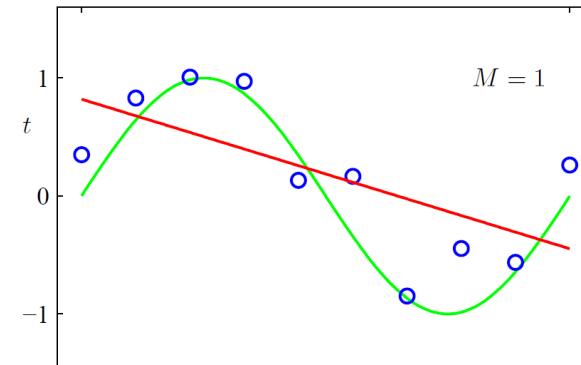
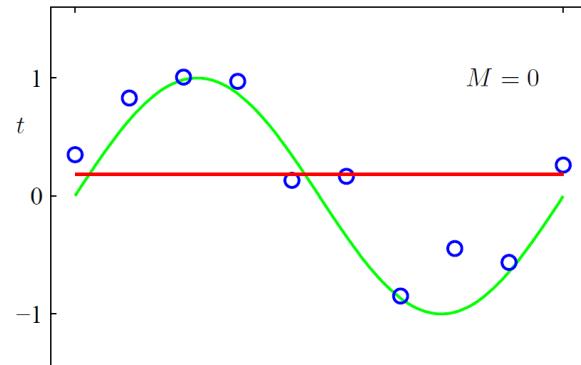


Nonlinear regression: example

- Training data are generated from a sin function + noise

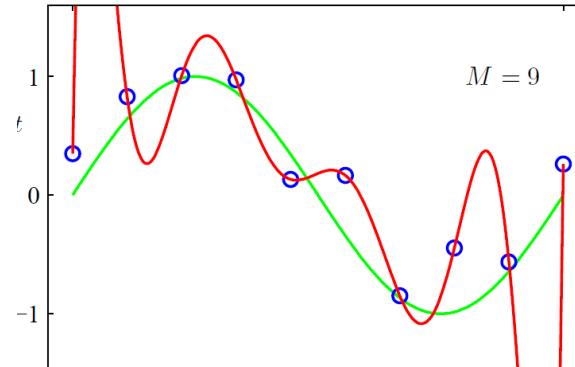
- Under-fitting

- h too simple



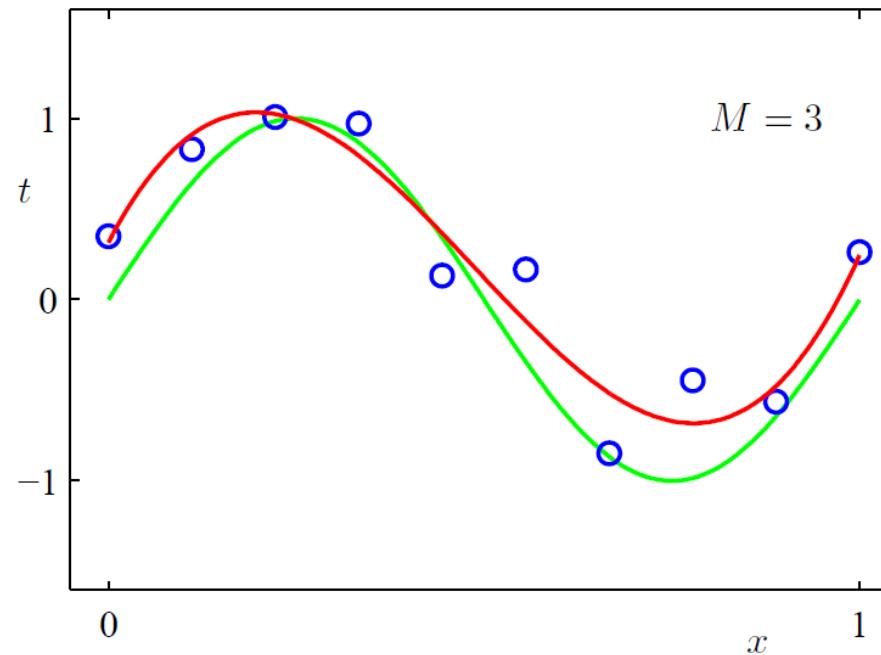
- Over-fitting

- h too complicated
training error = 0



Nonlinear regression: example

- Training data are generated from a sin function + noise
 - Just good enough



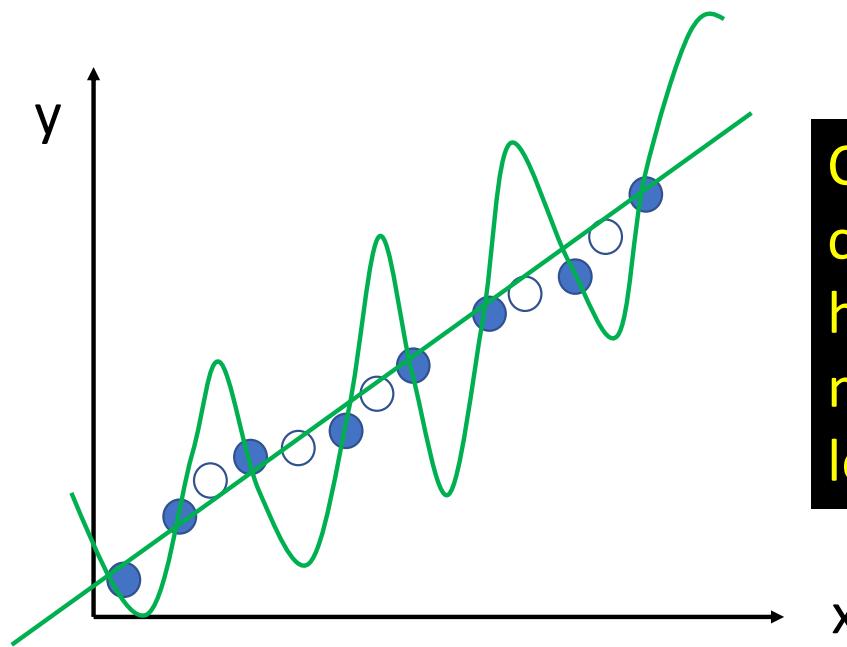
Training vs. testing

- Models learned from the training data should be applicable to test data

$$D_{train} = \{(x_1, y_1), \dots, (x_N, y_N)\}$$

$$D_{test} = \{(x_1, y_{N+1}), \dots, (x_N, y_{N+M})\}$$

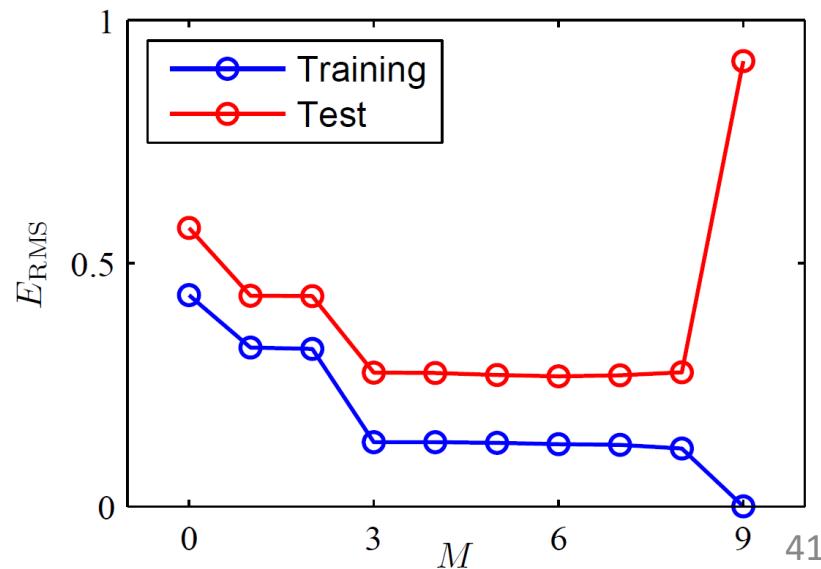
In training, we only see training data!



Choosing a more complicated hypothesis class does not necessarily lead to lower test errors!

Over-fitting vs. hypothesis class (brief)

- Over-fitting:
 - Small training error
 - Large test error (even larger than some other “simpler” models)
- How to quantify hypothesis class’ complexity?
 - KNN: smaller K , larger complexity
 - Polynomial regression: larger M , larger complexity
- How to choose?



Nonlinear regression: example

- 2nd-order polynomial regression for 2-D x :

$$\phi(x) = \begin{bmatrix} 1 \\ x[1] \\ x[2] \\ x[1]^2 \\ x[2]^2 \\ x[1]x[2] \end{bmatrix}$$

CSE 5523: Optimization



THE OHIO STATE UNIVERSITY

Hw-1

- Released on 9/3 midnight
- Due on 9/17 midnight
- Two questions in the programming set; five questions in the problem set
- Please use Piazza or come to the office hours for discussion
- Please start working on it ASAP

Next Tuesday (9/10)

- I will be traveling
- The class will be through Zoom, pre-recorded, or canceled.
 - Will announce it soon.
- Midterm
 - Will announce the date soon, by this weekend.

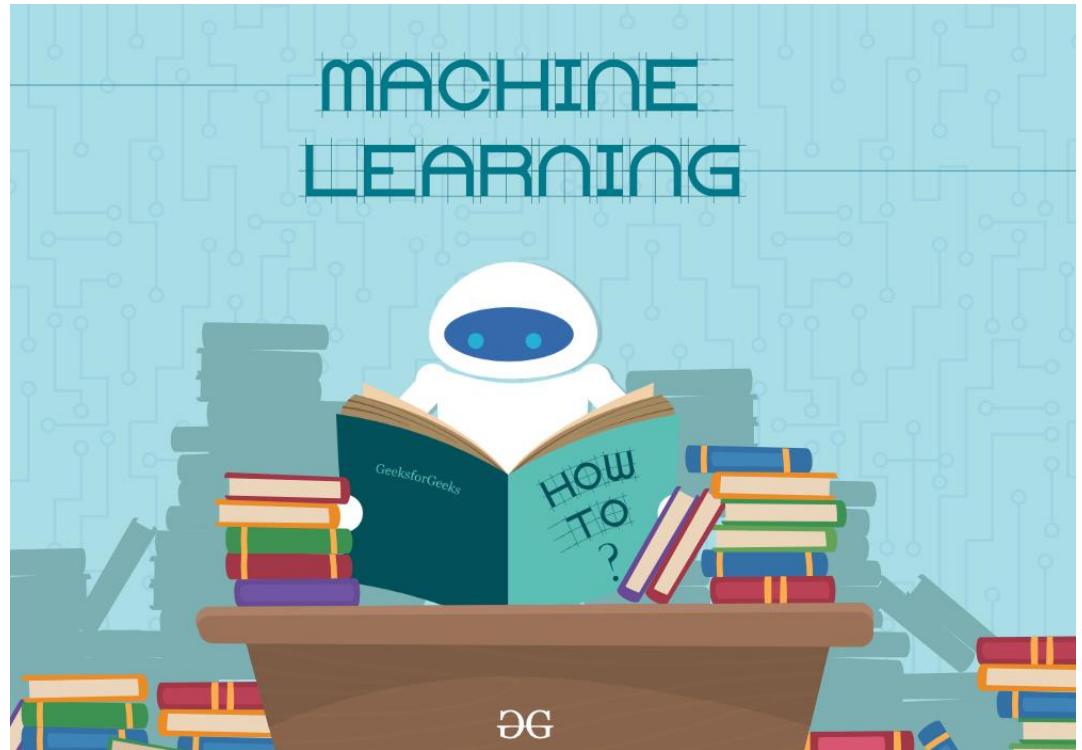
Today

Review

Nonlinear regression

Optimization

- Gradient descent
- Newton's method



Least mean square (LMS) solution

- $\tilde{\mathbf{w}}^{\text{LMS}} = \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \text{RSS}(\tilde{\mathbf{w}}) = \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_i [y_i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i]^2$

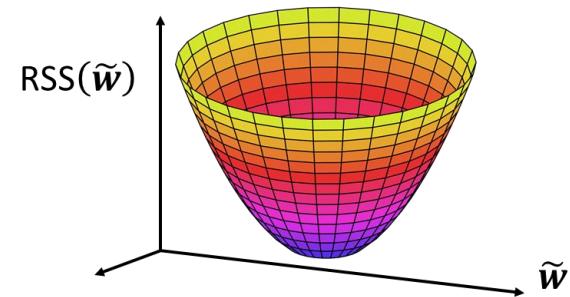
- Into a matrix representation

$$\text{RSS}(\tilde{\mathbf{w}}) = \sum_i (y_i - \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i)^2 = \tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - 2(\tilde{\mathbf{X}} \mathbf{y})^T \tilde{\mathbf{w}} + \text{constant}$$

- Normal equations

$$\frac{\partial \text{RSS}(\tilde{\mathbf{w}})}{\partial \tilde{\mathbf{w}}} = 2\tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - 2\tilde{\mathbf{X}} \mathbf{y} = \mathbf{0} \in \mathbb{R}^{(D+1)}$$

Closed-form $\tilde{\mathbf{w}}^{\text{LMS}} = (\tilde{\mathbf{X}} \tilde{\mathbf{X}}^T)^{-1} \tilde{\mathbf{X}} \mathbf{y}$



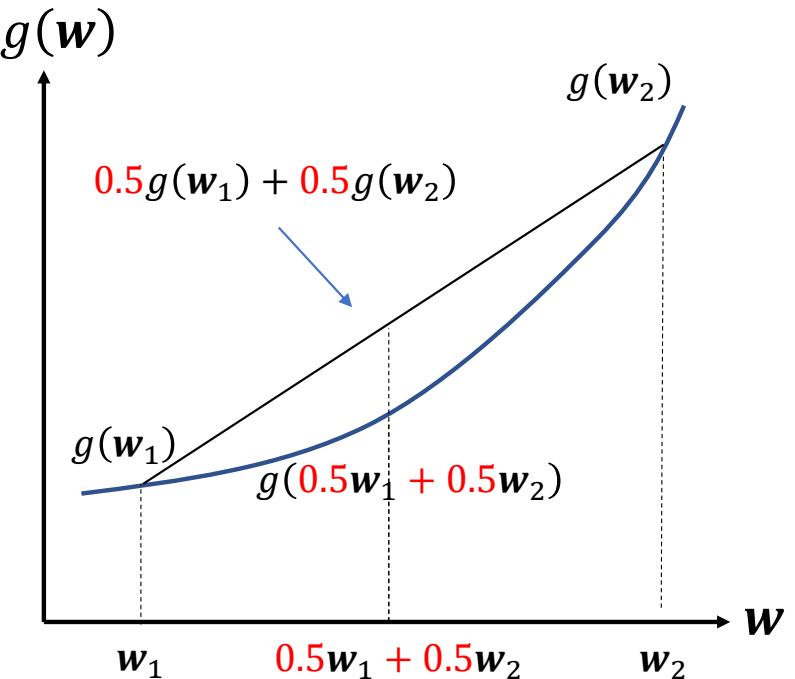
A little more math: convex functions

- A function $g(\mathbf{w})$ is convex iff $g(t\mathbf{w}_1 + (1 - t)\mathbf{w}_2) \leq tg(\mathbf{w}_1) + (1 - t)g(\mathbf{w}_2)$
 - $\forall \mathbf{w}_1, \mathbf{w}_2; \forall t \in [0, 1]$

- If $g(\mathbf{w})$ is twice differentiable, then

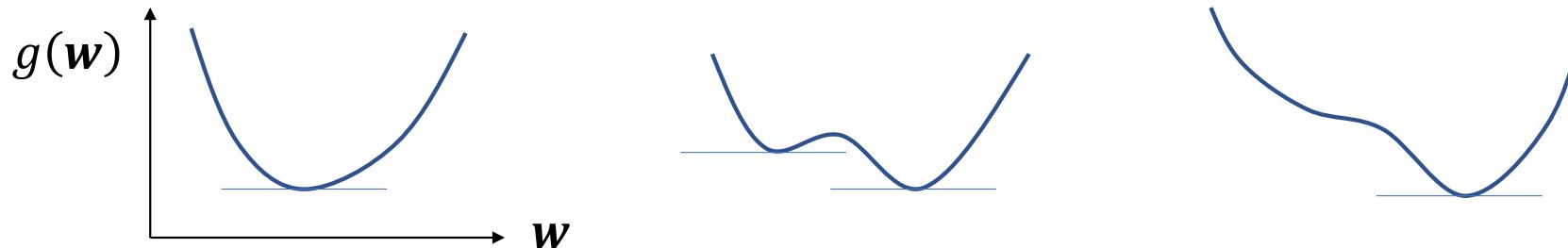
- $\mathbf{H}g(\mathbf{w}) = \begin{bmatrix} \frac{\partial^2 g(\mathbf{w})}{\partial w[1]\partial w[1]} & \dots & \frac{\partial^2 g(\mathbf{w})}{\partial w[1]\partial w[D]} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 g(\mathbf{w})}{\partial w[D]\partial w[1]} & \dots & \frac{\partial^2 g(\mathbf{w})}{\partial w[D]\partial w[D]} \end{bmatrix} \geq 0, \forall \mathbf{w}$

- Hessian matrix is positive semi-definite (PSD), $\forall \mathbf{w}$



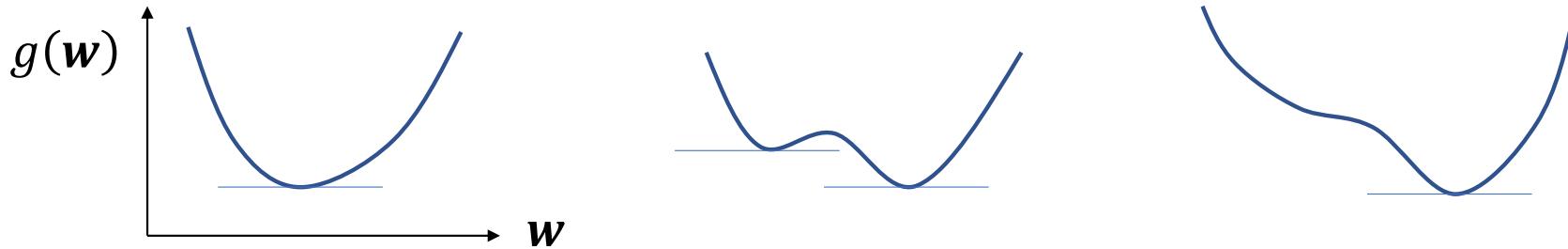
A little more math: convex functions

- Properties: Any local optimum is a global optimum
 - Gradients (partial derivatives) = 0 means local optimums



- Is $\text{RSS}(\tilde{\mathbf{w}})$ convex?
 - $\text{RSS}(\tilde{\mathbf{w}}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - 2(\tilde{\mathbf{X}}\mathbf{y})^T \tilde{\mathbf{w}} + \text{constant}$
 - $\nabla \text{RSS}(\tilde{\mathbf{w}}) = 2\tilde{\mathbf{X}} \tilde{\mathbf{X}}^T \tilde{\mathbf{w}} - 2\tilde{\mathbf{X}}\mathbf{y}$
 - $\mathbf{H} \text{RSS}(\tilde{\mathbf{w}}) = 2\tilde{\mathbf{X}} \tilde{\mathbf{X}}^T$
 - PSD because $\mathbf{v}^T (2\tilde{\mathbf{X}} \tilde{\mathbf{X}}^T) \mathbf{v} = 2(\tilde{\mathbf{X}}^T \mathbf{v})^T \tilde{\mathbf{X}}^T \mathbf{v} = 2\|\tilde{\mathbf{X}}^T \mathbf{v}\|_2^2 \geq 0, \forall \mathbf{v} \in \mathbb{R}^{D+1} \neq \mathbf{0}$

Convex functions



- When the function to be minimized is convex:
 - Any local optimum is a global optimum
 - However, it does not imply there exists closed-form solutions
 - Indeed, most of the (convex) optimization problems do not have closed-form solutions

Questions?



THE OHIO STATE UNIVERSITY

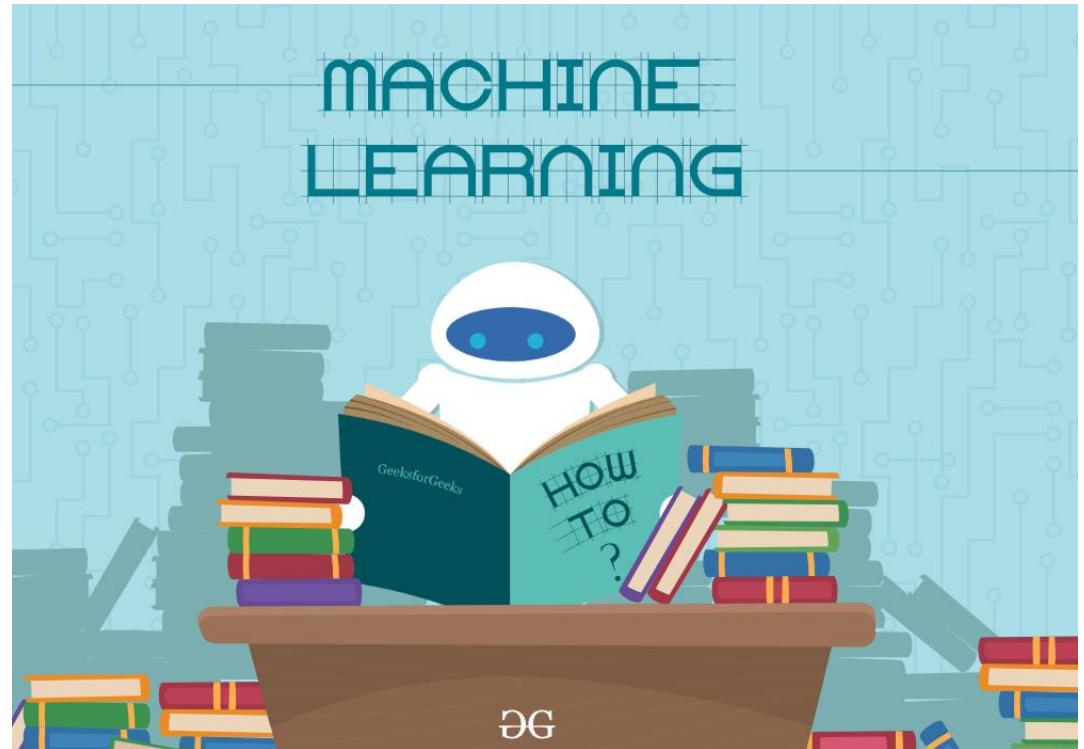
Today

Review

Nonlinear regression

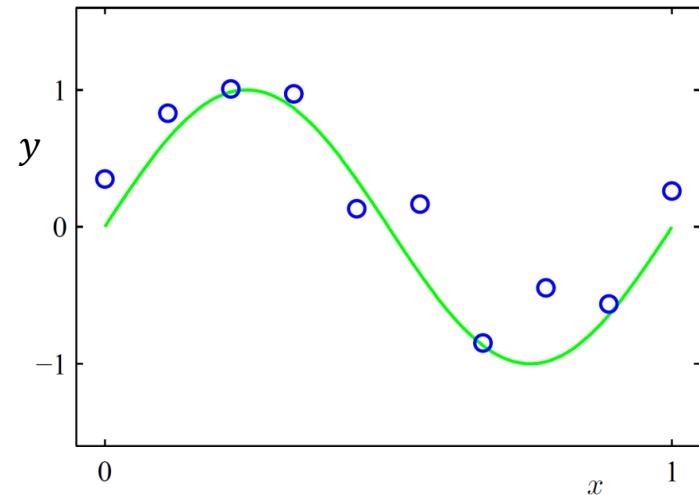
Optimization

- Gradient descent
- Newton's method



Nonlinear regression

- Definition: fit a **nonlinear curve** to data
 - $h(\mathbf{x})$ is linear iff $h(a\mathbf{x}_1 + b\mathbf{x}_2) = a \times h(\mathbf{x}_1) + b \times h(\mathbf{x}_2)$
 - $\forall \mathbf{x}_1, \mathbf{x}_2; \forall a, b \in \mathbb{R}$
 - Otherwise, it is nonlinear
 - Example: $h(\mathbf{x}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$ is linear
- Questions:
 - How to represent a nonlinear function?
 - What are the parameters?
 - How to estimate the parameters?



Nonlinear regression: representation

- Method 1: directly use nonlinear functions

- $h(x; w) = e^{wx}$
- $h(x; w) = \sin(wx)$
- ...

Explicitly indicate the parameters;
alternatively, $h_w(x)$

- Method 2: use nonlinear **feature transforms** or **mappings**

- $\phi(x): x \in \mathbb{R}^D \mapsto z \in \mathbb{R}^M$
- M is the dimensionality of the transformed feature z (or $\phi(x)$); can be larger, smaller than, or the same as D
- Notation $z[m] = \phi_m(x)$
- ϕ_m is called a nonlinear basis function
- $\phi(x)$ can be used for other ML algorithms (e.g., nearest neighbor) as well

Nonlinear regression: representation

- Examples of $\phi_m(x)$:
 - $z[m] = (x[m])^2$
 - $z[m] = \sin(x[m])$
 - $z[m] = \log(x[m])$
 - $z[m] = e^{x[m]}$
 - When $x \in \mathbb{R}$
 - $z[m] = (x)^{m-1}$
 - $z \in \mathbb{R}^M$ can already encode the dummy variable 1

Nonlinear regression: representation

- Previously we use $h(\mathbf{x}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$
- Now we use $h(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) = \sum_m w[m] \phi_m(\mathbf{x})$, where $\mathbf{w} \in \mathbb{R}^M$
 - In nonlinear regression, the number of parameters is usually different from D or D+1
- $\text{RSS}(\mathbf{w}) = \sum_i [y_i - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i)]^2$
- $\mathbf{w}^{\text{LMS}} = \underset{\mathbf{w}}{\operatorname{argmin}} \text{RSS}(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_i [y_i - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i)]^2$

Nonlinear regression: parameter estimation

- $\mathbf{w}^{\text{LMS}} = \underset{\mathbf{w}}{\operatorname{argmin}} \text{RSS}(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_i [y_i - \mathbf{w}^T \boldsymbol{\phi}(x_i)]^2$

$$\Phi \in \mathbb{R}^{M \times N} = [\boldsymbol{\phi}(x_1) \quad \dots \quad \boldsymbol{\phi}(x_N)] = \begin{bmatrix} | & & | \\ \boldsymbol{\phi}(x_1) & \dots & \boldsymbol{\phi}(x_N) \\ | & & | \end{bmatrix}$$

- Closed form solution \mathbf{w}^{LMS} : $(\tilde{\Phi} \tilde{\Phi}^T)^{-1} \tilde{\Phi} \mathbf{y}$

- Question:

- How many data are needed?
 - Is $\text{RSS}(\mathbf{w})$ still convex?

Questions?



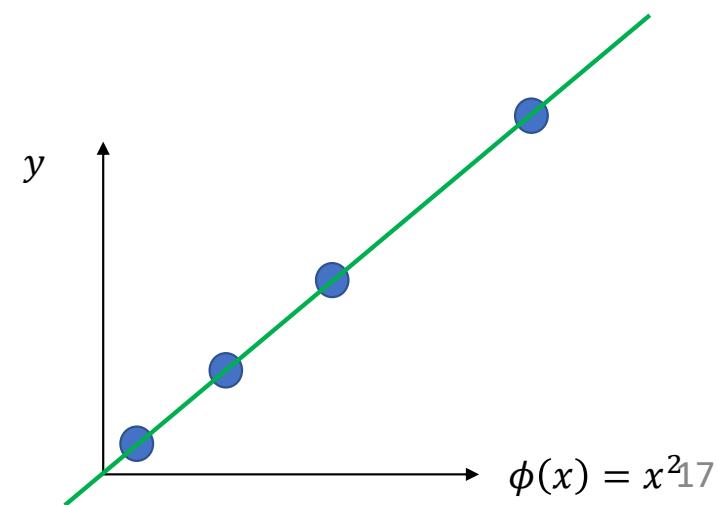
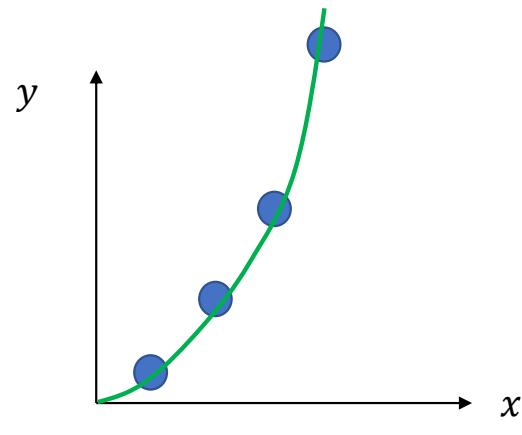
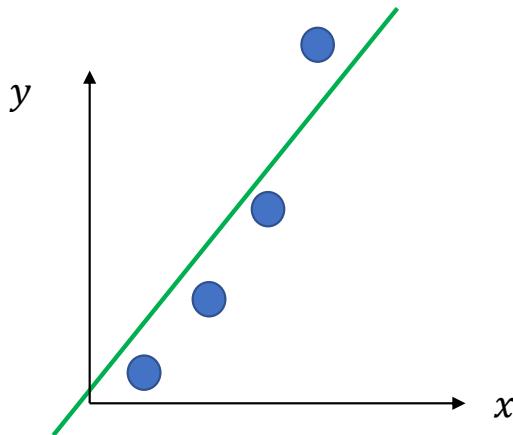
THE OHIO STATE UNIVERSITY

Nonlinear regression: example

- Polynomial regression for 1-D x , i.e., $h(x) = \sum_{m=0}^M w[m]x^m = \mathbf{w}^T \boldsymbol{\phi}(x)$

$$\boldsymbol{\phi}(x) = \begin{bmatrix} 1 \\ x \\ \vdots \\ x^M \end{bmatrix}$$

- Special case: $D_{tr} = \{(x_n \in \mathbb{R}, y_n = x_n^2)\}_{n=1}^N$

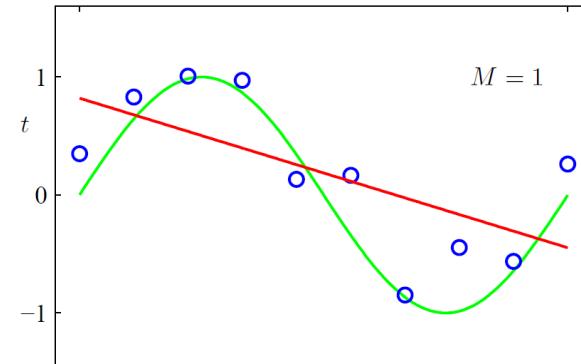
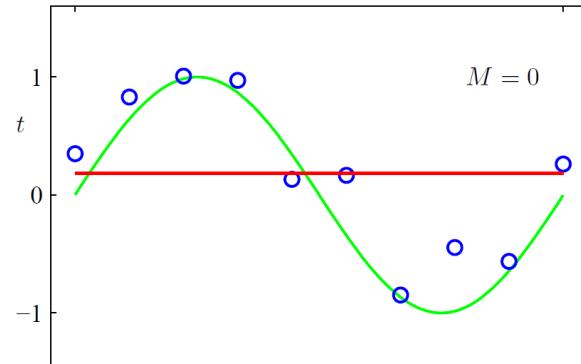


Nonlinear regression: example

- Training data are generated from a sin function + noise

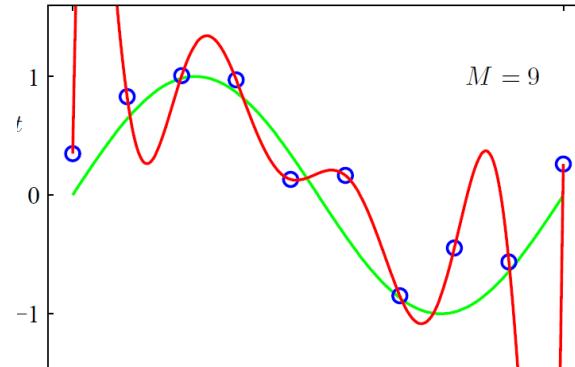
- Under-fitting

- h too simple



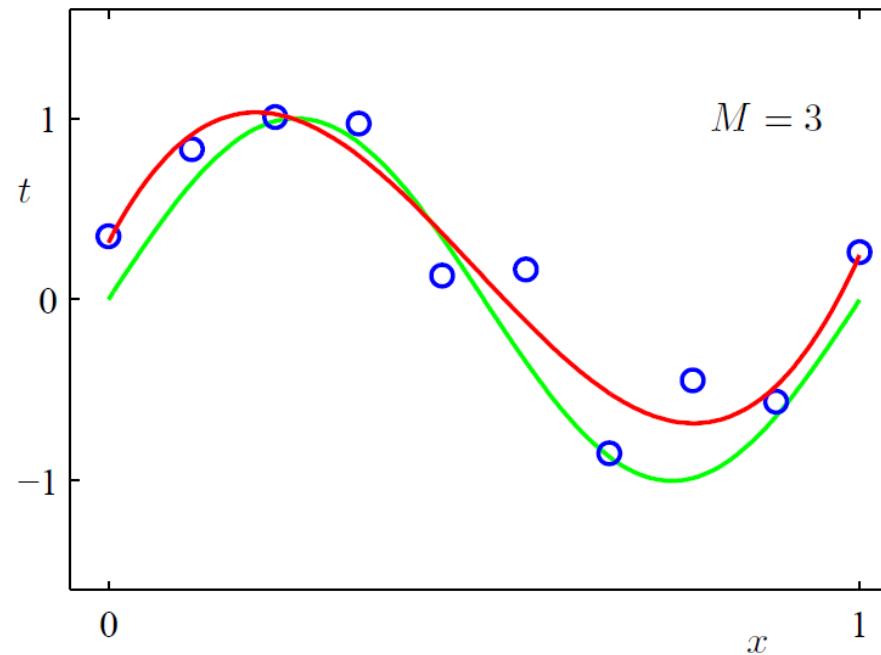
- Over-fitting

- h too complicated
training error = 0



Nonlinear regression: example

- Training data are generated from a sin function + noise
 - Just good enough



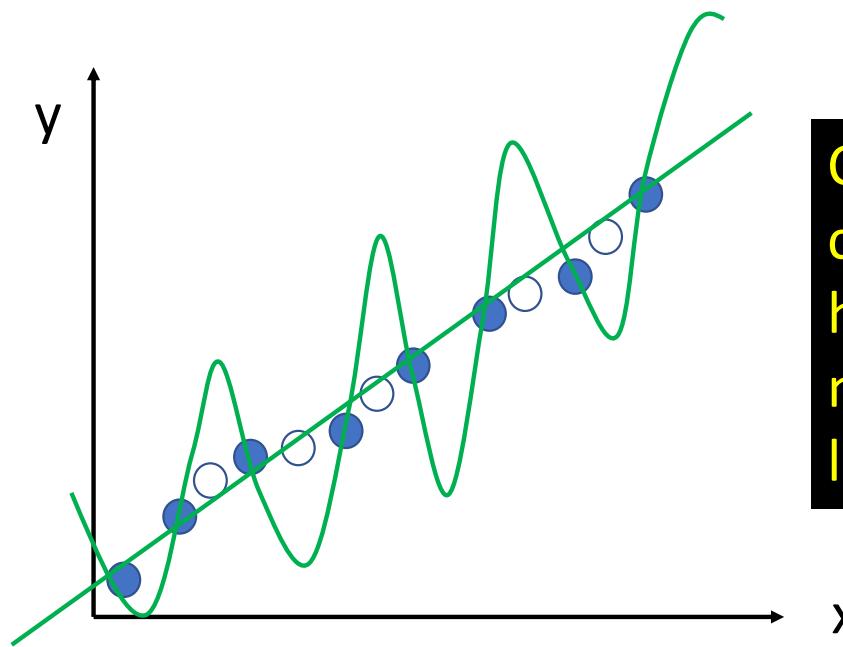
Training vs. testing

- Models learned from the training data should be applicable to test data

$$D_{train} = \{(x_1, y_1), \dots, (x_N, y_N)\}$$

$$D_{test} = \{(x_1, y_{N+1}), \dots, (x_N, y_{N+M})\}$$

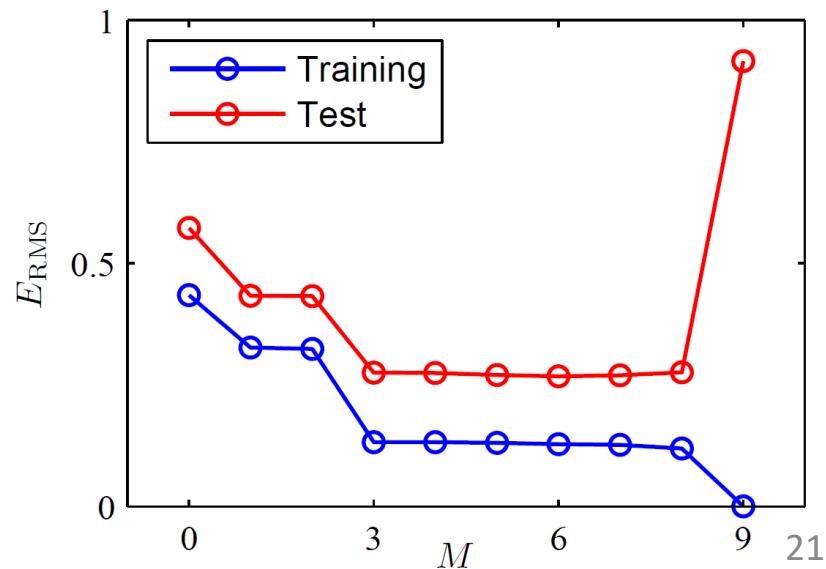
In training, we only see training data!



Choosing a more complicated hypothesis class does not necessarily lead to lower test errors!

Over-fitting vs. hypothesis class (brief)

- Over-fitting:
 - Small training error
 - Large test error (even larger than some other “simpler” models)
- How to quantify hypothesis class’ complexity?
 - KNN: smaller K , larger complexity
 - Polynomial regression: larger M , larger complexity
- How to choose?



Nonlinear regression: example

- 2nd-order polynomial regression for 2-D x :

$$\phi(x) = \begin{bmatrix} 1 \\ x[1] \\ x[2] \\ x[1]^2 \\ x[2]^2 \\ x[1]x[2] \end{bmatrix}$$

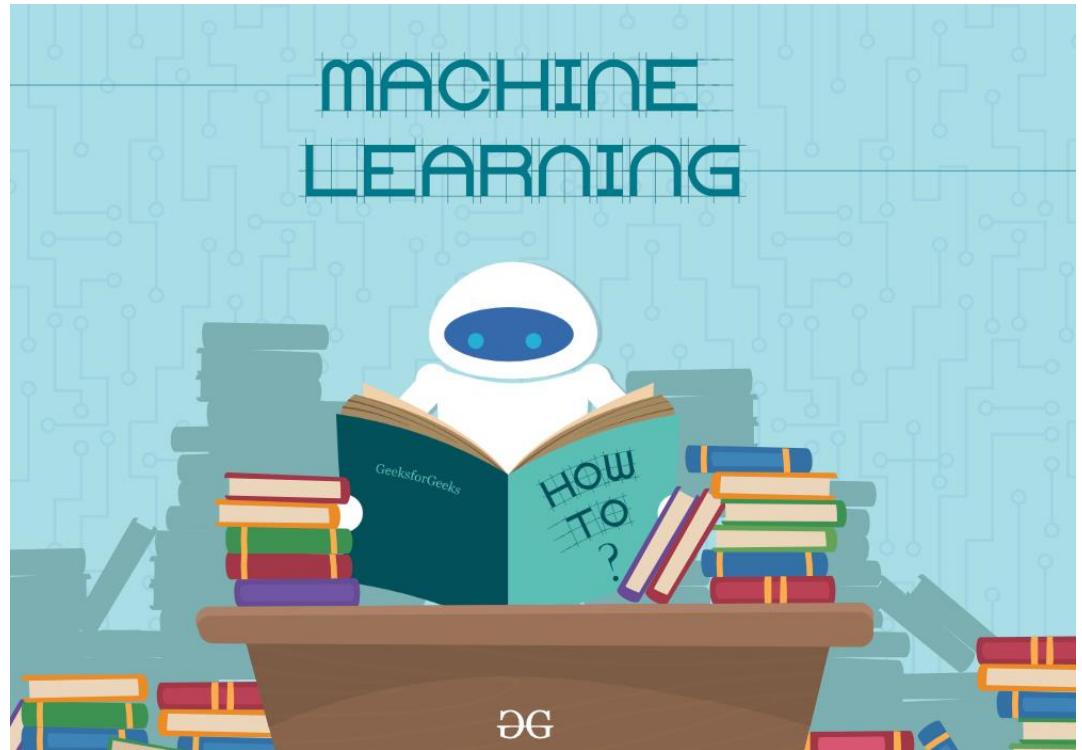
Today

Review

Nonlinear regression

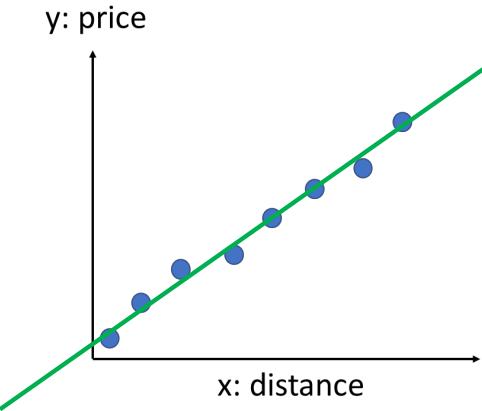
Optimization

- Gradient descent
- Newton's method



In finding the solution (estimating parameters) ...

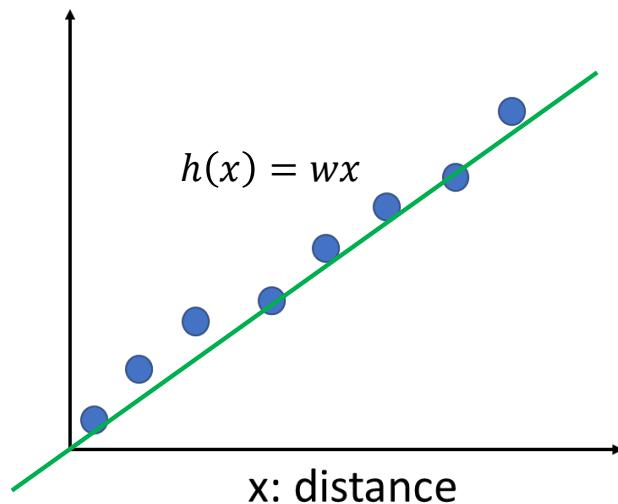
- The training data are fixed
- The variables are the parameters \mathbf{w} to estimate
- A “machine learning” algorithm finds \mathbf{w} to minimize the training loss
 - **Question: how to do this?**
- For the “users”, we care about given a future x , what the prediction will be



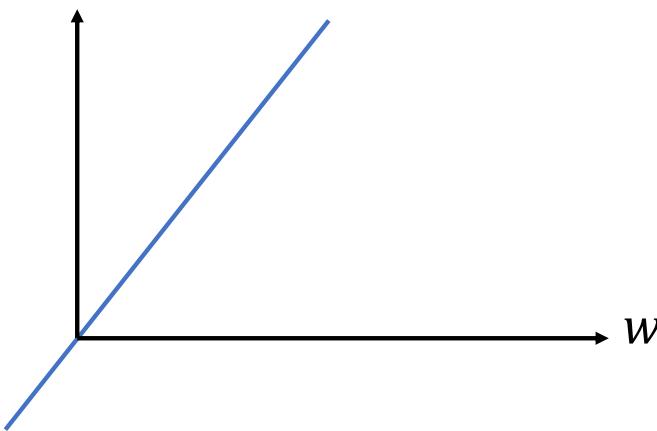
$$E(\mathbf{w}) = \sum_i [y_i - h(x_i; \mathbf{w})]^2$$

Different figures

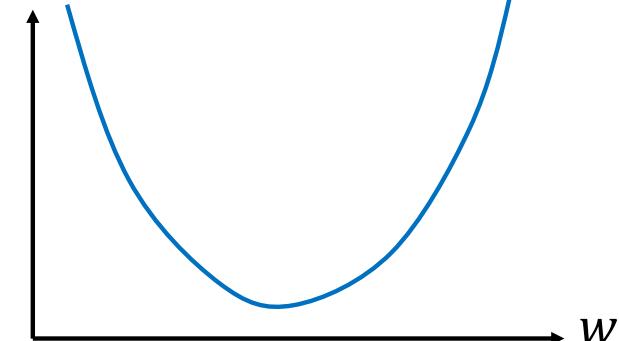
y: price



$h(x) = wx$ for a certain x



$$E(w) = \sum_i (wx_i - y_i)^2$$



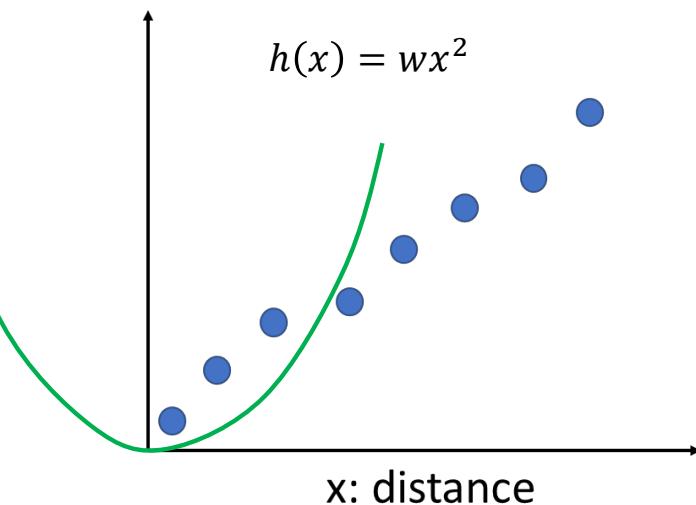
Convex optimizations require the error function, not the prediction function, to be convex!

Linear w.r.t. w !

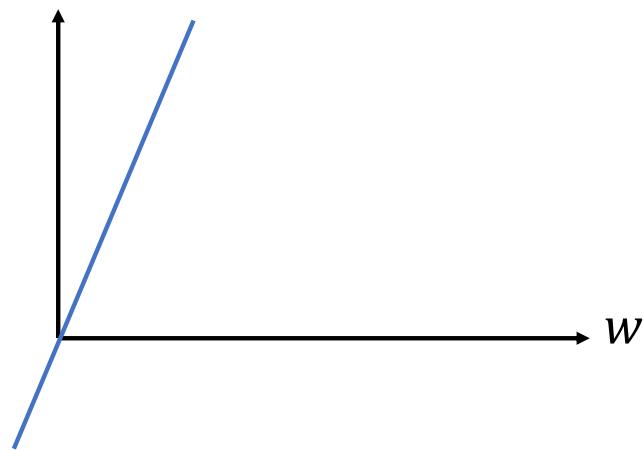
Convex

Different figures

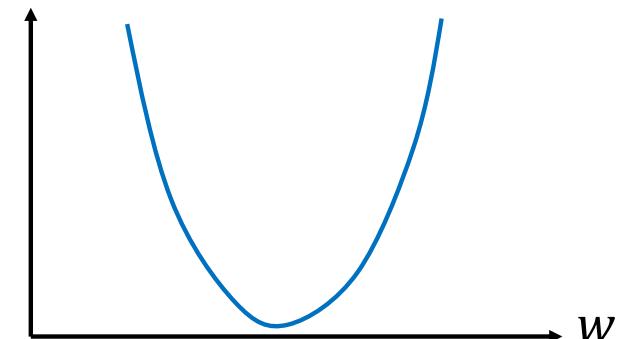
y: price



$h(x) = wx^2$ for a certain x



$$E(w) = \sum_i (wx_i^2 - y_i)^2$$



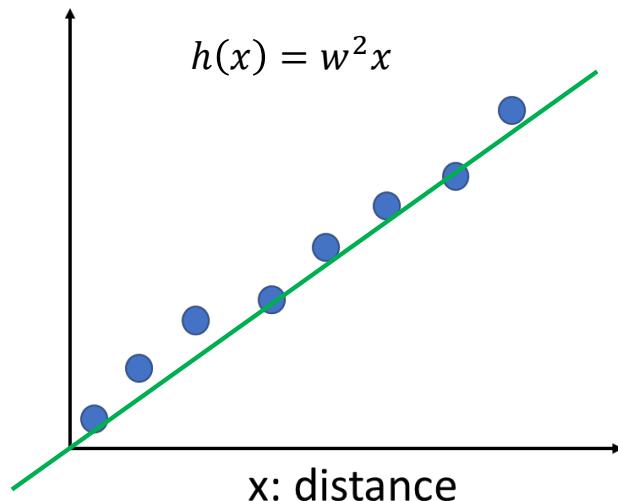
Nonlinear regression since
 $h(x) = wx^2$ is not linear w.r.t. x .

Still linear w.r.t. w !

Convex

Different figures

y: price



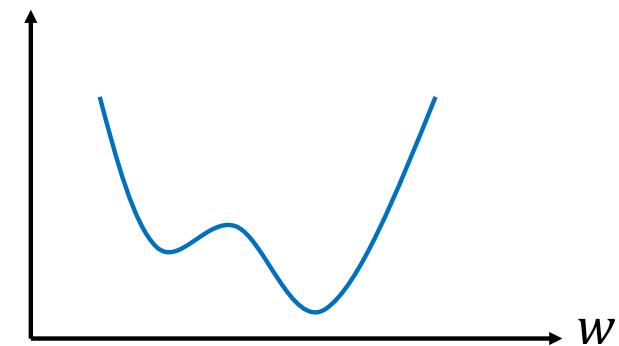
Linear regression since
 $h(x) = w^2x$ is linear w.r.t. x .

$h(x) = w^2x$ for a certain x



Not linear w.r.t. w !

$$E(w) = \sum_i (w^2 x_i - y_i)^2$$



Maybe not convex

Questions?

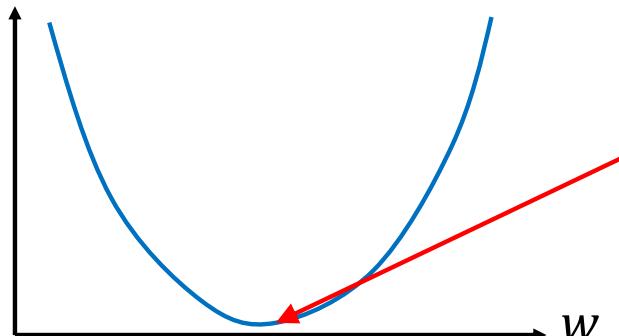


THE OHIO STATE UNIVERSITY

When $h(x; w)$ is nonlinear w.r.t. w

- $\text{RSS}(w) = \sum_i (y_i - h(x_i; w))^2$ is no longer parabolic w.r.t. the parameters w
 - **May be nonconvex:** local minimums are not always global minimums
 - **May have no closed-form solution:** need numerical approximation

$$E(w) = \sum_i (wx_i - y_i)^2$$

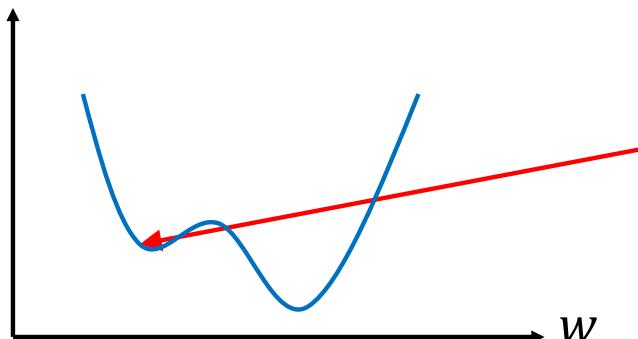


- ❖ $\frac{\partial E}{\partial w} = 0$ means global minimums
- ❖ Closed-forms to find w for $\frac{\partial E}{\partial w} = 0$

When $h(x; w)$ is nonlinear w.r.t. w

- $\text{RSS}(w) = \sum_i (y_i - h(x_i; w))^2$ is no longer parabolic w.r.t. the parameters w
 - **May be nonconvex:** local minimums are not always global minimums
 - **May have no closed-form solution:** need numerical approximation

$$E(w) = \sum_i (h(x_i; w) - y_i)^2$$



- ❖ $\frac{\partial E}{\partial w} = 0$ does not mean global minimums
- ❖ Maybe no closed-forms to find w for $\frac{\partial E}{\partial w} = 0$

When other loss functions are used

- The training loss $E(\mathbf{w})$ may not be parabolic
 - May have no closed-form solution for $\frac{\partial E}{\partial \mathbf{w}} = 0$
- The training loss $E(\mathbf{w})$ may not be convex
 - $\frac{\partial E}{\partial \mathbf{w}} = 0$ does not mean global minimums
- How are we going to do parameter estimation?

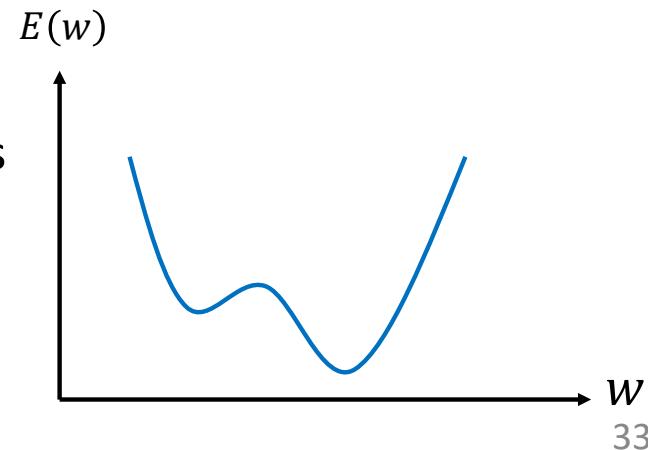
Questions?



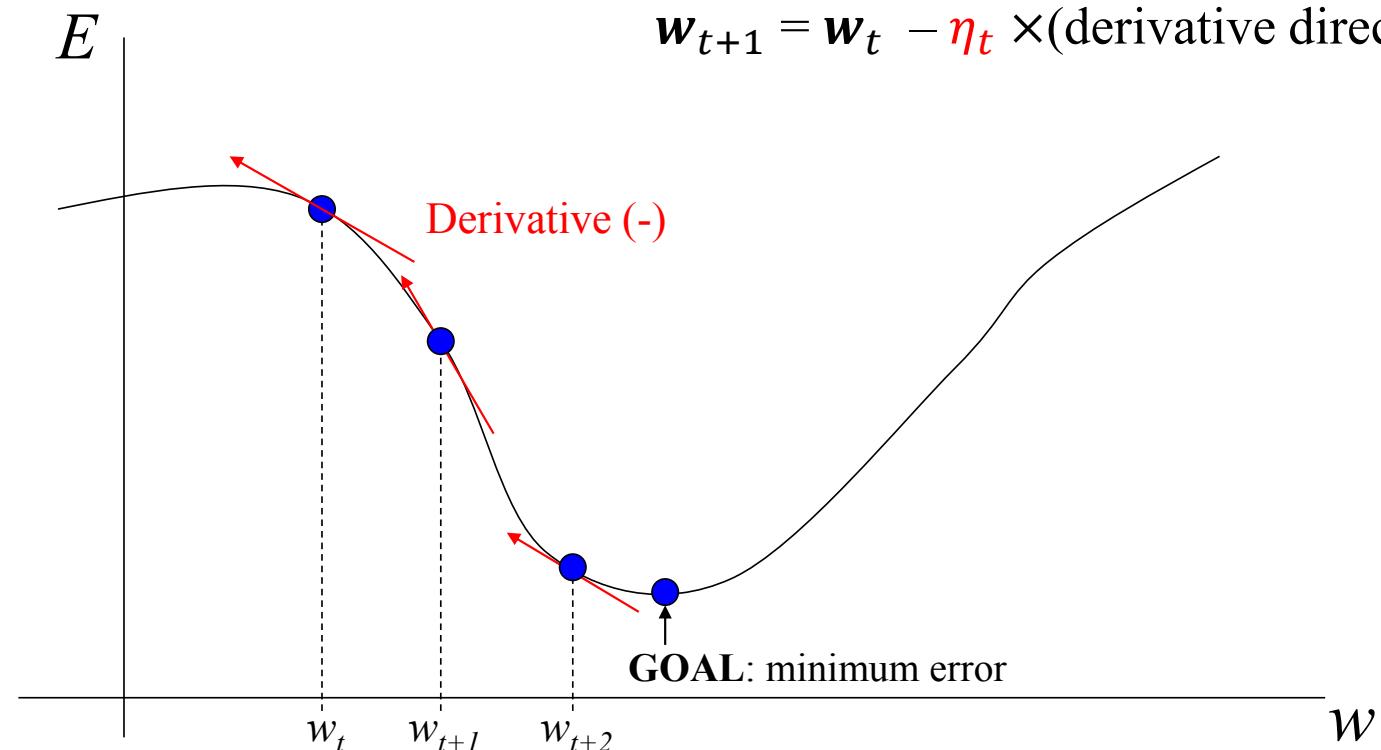
THE OHIO STATE UNIVERSITY

Gradient descent (GD)

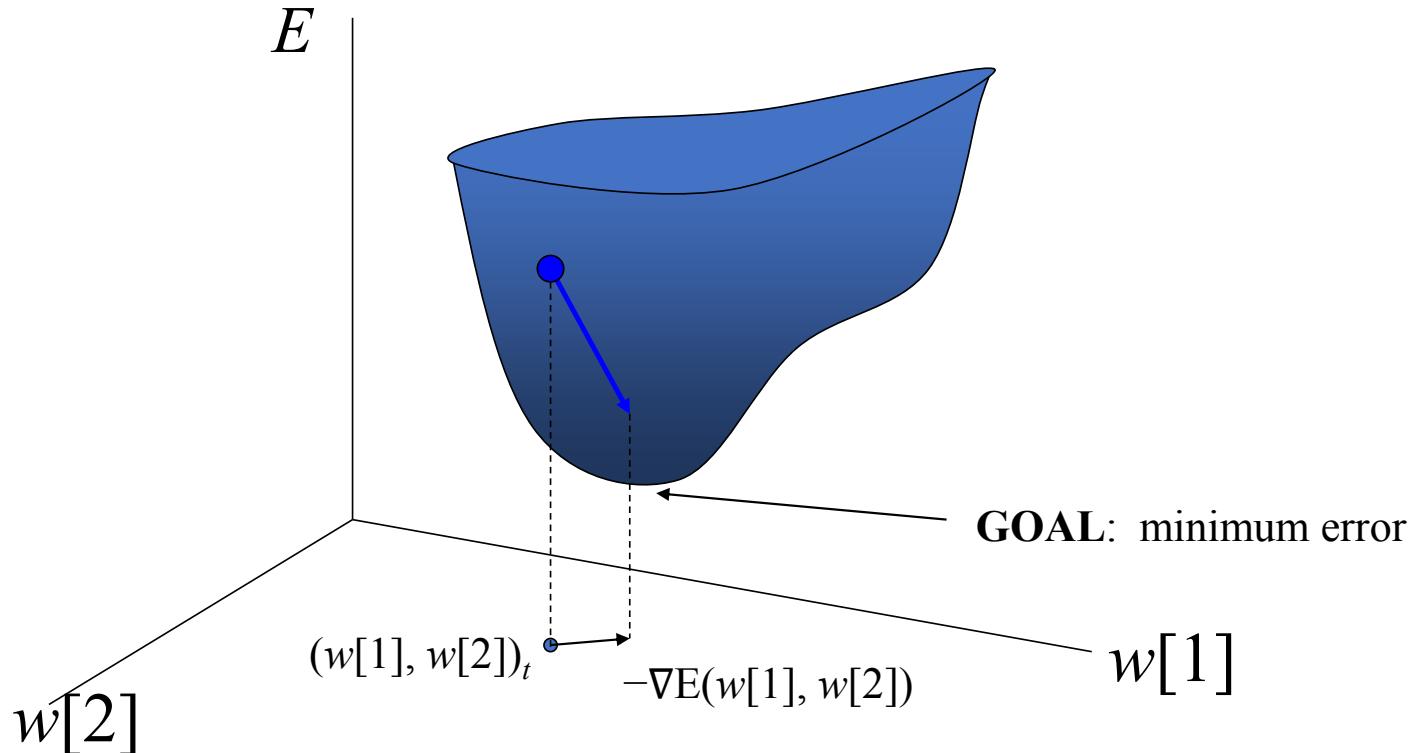
- A simple (and fast?) method for optimization
 - $h(\mathbf{x}; \mathbf{w})$ linear or nonlinear w.r.t. parameters \mathbf{w}
 - Sum or square losses or other losses
- A general method to find local minimum for a multi-parameter error function
 - Search through “parameter space” to find the minimum error
- Main idea
 - Evaluate the error function for current guess of parameters
 - Determine change in parameters $\Delta\mathbf{w}$ that decreases error
 - From gradient of error function
 - Update parameters with this new change
 - Repeat process until converge (or hit max iterations)



Problem visualization



Problem visualization



Local approximation (let w be M -dim)

- Taylor series expansion

$$E(w) = E(w') + \nabla E(w')^T (w - w') + \frac{1}{2} (w - w')^T H_E(w') (w - w') + H.O.T.$$

Constant term
Linear (first-order) term
Quadratic (second-order) term

Variable to solve
Current guess
Variation “around” w'

- Rewrite $w = (w' + \Delta w)$

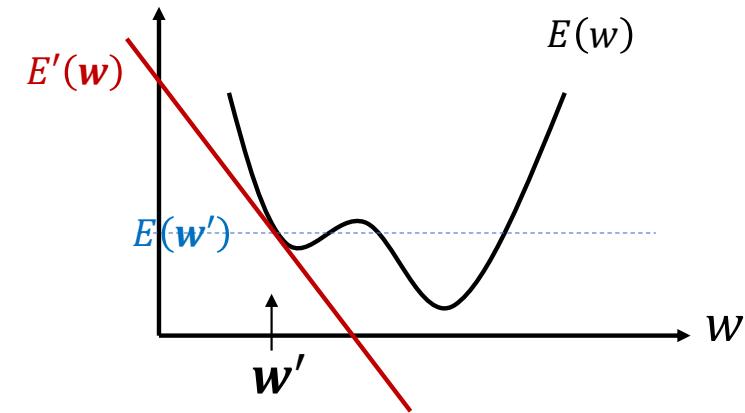
$$E(w' + \Delta w) = E(w') + \nabla E(w')^T \Delta w + \frac{1}{2} \Delta w^T H_E(w') \Delta w + H.O.T.$$

Δw : the movement from w'

Local approximation (let w be M -dim)

- Visualization

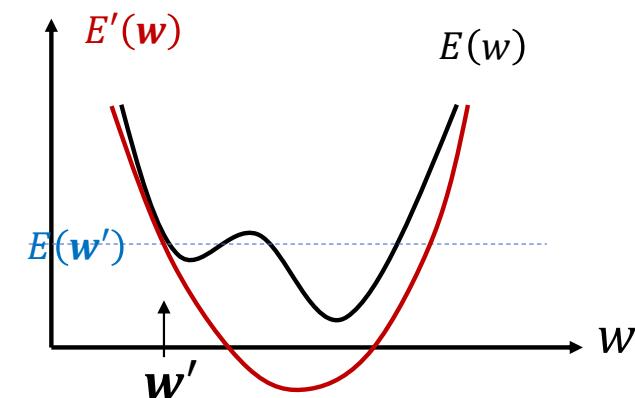
- First order: $E'(w) = E(w') + \nabla E(w')^T (w - w')$



- Second order: $E'(w) = E(w') + \nabla E(w')^T (w - w')$

$$+ \frac{1}{2} (w - w')^T \mathbf{H}_E(w') (w - w')$$

- The approximation $E'(w)$ is quadratic!
 - If $E(w)$ is strictly convex, its Hessian matrix $\mathbf{H}_E(w') > 0$: closed forms to minimize the approximation $E'(w)$



Why does gradient descent work?

- Recall the Taylor expansion

$$E(\mathbf{w}) = E(\mathbf{w}') + \nabla E(\mathbf{w}')^T (\mathbf{w} - \mathbf{w}') + \text{H.O.T.} = E(\mathbf{w}') + \nabla E(\mathbf{w}')^T \Delta\mathbf{w} + \text{H.O.T.}$$

- Approximate $E(\mathbf{w})$ by the linear terms: $E'(\mathbf{w}' + \Delta\mathbf{w}) = E(\mathbf{w}') + \nabla E(\mathbf{w}')^T \Delta\mathbf{w}$
- Which direction, i.e., $\Delta\mathbf{w}$ with $\|\Delta\mathbf{w}\|_2^2 = 1$, can minimize $E'(\mathbf{w}' + \Delta\mathbf{w})$?
 - Answer: **direction of** $\Delta\mathbf{w} = -\frac{\nabla E(\mathbf{w}')}{\|\nabla E(\mathbf{w}')\|_2}$
 - Interpretation: move in the direction of $-\nabla E(\mathbf{w}')$ is the most efficient for minimization
 - Only works for **small $\Delta\mathbf{w}$ (so a scale η)** due to the nature of linear Taylor expansion

The GD Algorithm

Initialize parameters:

$$\hat{\mathbf{w}}^{(0)} = [w[1], w[2], \dots, w[M]]^T$$

Error function:

$$E(\mathbf{w})$$

Compute gradients:

$$\nabla E(\hat{\mathbf{w}}^{(t)}) = \begin{bmatrix} \frac{\partial E(\hat{\mathbf{w}}^{(t)})}{\partial w[1]} \\ \vdots \\ \frac{\partial E(\hat{\mathbf{w}}^{(t)})}{\partial w[M]} \end{bmatrix}$$

Update parameters: $\hat{\mathbf{w}}^{(t+1)} = \hat{\mathbf{w}}^{(t)} - \eta^{(t)} \cdot \nabla E(\hat{\mathbf{w}}^{(t)})$

[Repeat until converges]

Learning rate

2-order Taylor expansion approximation

- Taylor series expansion ($\Delta\mathbf{w} = \mathbf{w} - \mathbf{w}'$)

$$E(\mathbf{w}' + \Delta\mathbf{w}) = E(\mathbf{w}') + \nabla E(\mathbf{w}')^T \Delta\mathbf{w} + \frac{1}{2} \Delta\mathbf{w}^T \mathbf{H}_E(\mathbf{w}') \Delta\mathbf{w} + \text{H.O.T.}$$

- Hessian matrix $\mathbf{H}_E(\mathbf{w}')$

$$\mathbf{H}_E(\mathbf{w}') = \begin{bmatrix} \frac{\partial^2 E(\mathbf{w}')}{\partial w[1]\partial w[1]} & \cdots & \frac{\partial^2 E(\mathbf{w}')}{\partial w[1]\partial w[M]} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 E(\mathbf{w}')}{\partial w[M]\partial w[1]} & \cdots & \frac{\partial^2 E(\mathbf{w}')}{\partial w[M]\partial w[M]} \end{bmatrix}$$

- Newton's method

- For any kind of “convex” loss $E(\mathbf{w})$
- Quadratic approximation of $E(\mathbf{w})$ at \mathbf{w}' : ignore H.O.T.
- Minimize the quadratic approximation to get $\Delta\mathbf{w} = -\mathbf{H}_E^{-1}(\mathbf{w}') \nabla E(\mathbf{w}')$
- Update \mathbf{w}' by $\mathbf{w}' + \eta^{(t)} \Delta\mathbf{w}$

- Pros: faster convergence; Cons: $\mathbf{H}_E(\mathbf{w}')$ and $\mathbf{H}_E^{-1}(\mathbf{w}')$ are time consuming

The Newton's method

Initialize parameters:

$$\hat{\mathbf{w}}^{(0)} = [w[1], w[2], \dots, w[M]]^T$$

Error function:

$$E(\mathbf{w})$$

Compute gradients and Hessian:

$$\nabla E(\hat{\mathbf{w}}^{(t)}) = \begin{bmatrix} \frac{\partial E(\hat{\mathbf{w}}^{(t)})}{\partial w[1]} \\ \vdots \\ \frac{\partial E(\hat{\mathbf{w}}^{(t)})}{\partial w[M]} \end{bmatrix}, \quad \mathbf{H}_E(\hat{\mathbf{w}}^{(t)}) = \begin{bmatrix} \frac{\partial^2 E(\hat{\mathbf{w}}^{(t)})}{\partial w[1] \partial w[1]} & \cdots & \frac{\partial^2 E(\hat{\mathbf{w}}^{(t)})}{\partial w[1] \partial w[M]} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 E(\hat{\mathbf{w}}^{(t)})}{\partial w[M] \partial w[1]} & \cdots & \frac{\partial^2 E(\hat{\mathbf{w}}^{(t)})}{\partial w[M] \partial w[M]} \end{bmatrix}$$

Update
parameters:

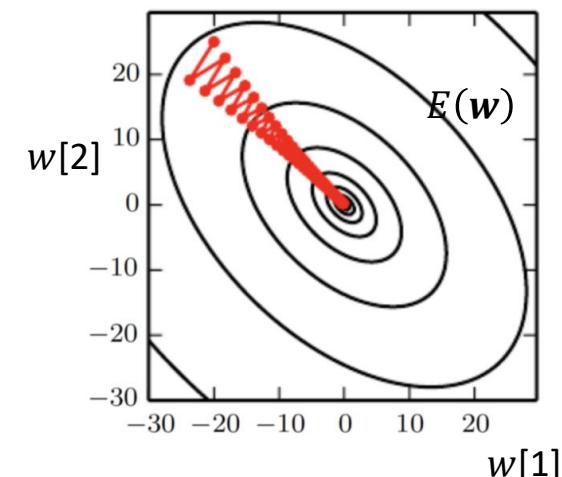
$$\hat{\mathbf{w}}^{(t+1)} = \hat{\mathbf{w}}^{(t)} - \eta^{(t)} \mathbf{H}_E^{-1}(\mathbf{w}') \nabla E(\mathbf{w}')$$

[Repeat until converges]

Learning rate

Summary

- Pros of GD:
 - No expensive matrix operations (i.e., inversion)
 - Non-RSS error functions (if gradient is well formed)
 - Easy to parallelize: $E(\mathbf{w}) = \sum_i \ell(h(\mathbf{x}_i; \mathbf{w}), y_i) \Rightarrow \nabla E(\mathbf{w}) = \sum_i \nabla \ell(h(\mathbf{x}_i; \mathbf{w}), y_i)$
- Cons of GD:
 - Can be misled by local minima and stuck at larger-error valleys
 - Can get stuck on a flat plain/flat region of the error function
 - Can overshoot back and forth (oscillate from side to side)
- Helpers
 - Use multiple random-restarts with random initial parameters
 - Methods exist for selecting appropriate (best) learning rate at each iteration



Summary

- Pros of Newton's method:
 - Non-RSS error functions (if gradient and Hessian are well formed)
 - If \mathbf{w}' is well-chosen, it can converge faster, since 2nd-order Taylor has more information.
- Cons of Newton's method:
 - Hessian and its inverse are time consuming
 - Divergence often happens if the function is flat or almost flat with respect to some dimension. In that case the second derivatives are close to zero, and their inverse becomes very large, resulting in gigantic steps.
- Helpers
 - Compute only the diagonal Hessian: $\Delta\mathbf{w} = -H^{-1}(\mathbf{w}')\nabla E(\mathbf{w}')$
 - GD and then Newton's method

Question

- How many steps do you need to optimize $E(w) = \sum_i (wx_i - y_i)^2$?
 - GD?
 - Newton's method?

CSE 5523: Perceptron



THE OHIO STATE UNIVERSITY

Hw-1

- Due on 9/17 midnight
- Two questions in the programming set; five questions in the problem set
- Please use Piazza or come to the office hours for discussion
- Please start working on the homework ASAP

Other reminders

- Midterm
 - 10/17 in class
- Office hour this Thursday
 - 6:00 pm – 7:30 pm
- Piazza
 - Still many of you not in Piazza
 - Please find the link and access code in Carmen Syllabus or Carmen Announcement

Today

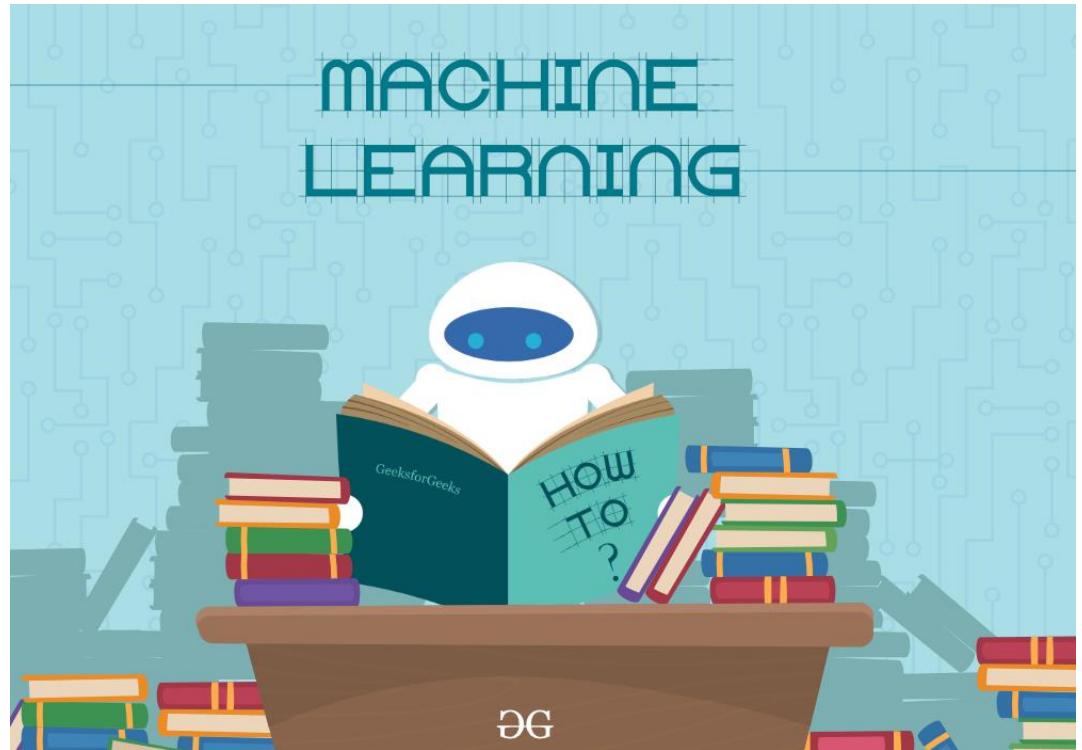
Optimization

- Gradient descent (continued)
- Newton's method

Linear classifier

- Definition
- First algorithm: perceptron
- Reading: pocket algorithm

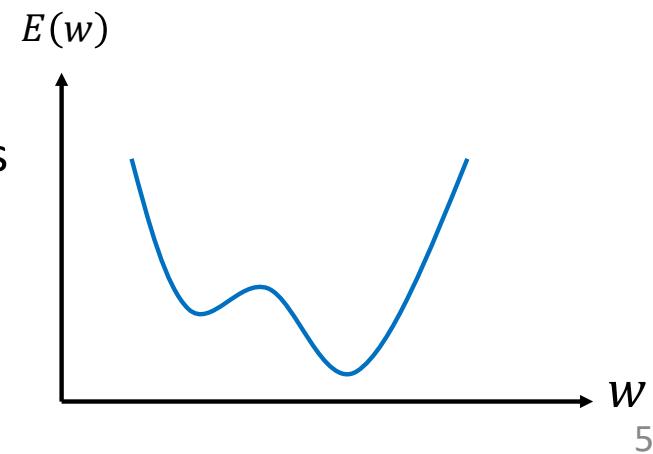
Probability refresher (self reading)



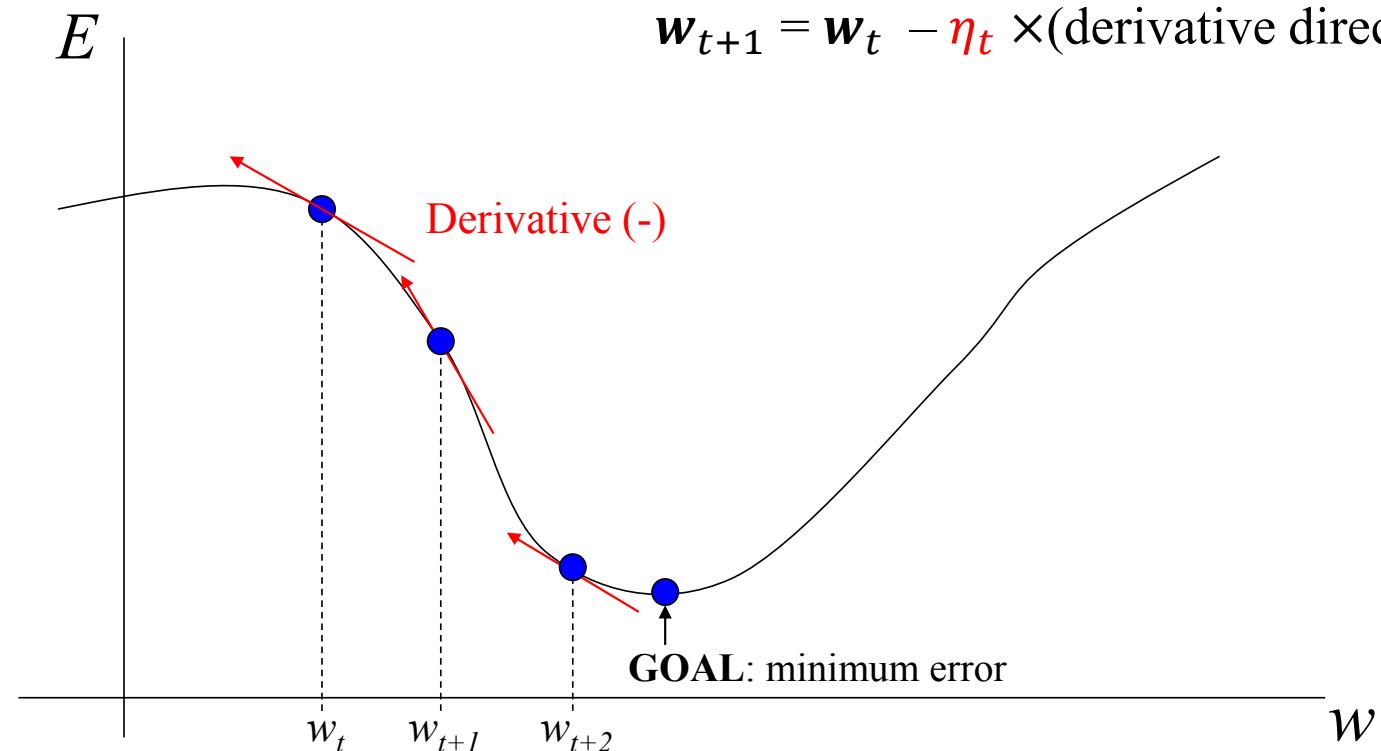
Gradient descent (GD)

- A simple (and fast?) method for optimization
 - $h(\mathbf{x}; \mathbf{w})$ linear or nonlinear w.r.t. parameters \mathbf{w}
 - Sum or square losses or other losses
- A general method to find local minimum for a multi-parameter error function
 - Search through “parameter space” to find the minimum error
- Main idea
 - Evaluate the error function for current guess of parameters
 - Determine change in parameters $\Delta\mathbf{w}$ that decreases error
 - From gradient of error function
 - Update parameters with this new change
 - Repeat process until converge (or hit max iterations)

$$E(\mathbf{w}) = \sum_i \ell(h(\mathbf{x}_i; \mathbf{w}), y_i)$$



Problem visualization



Local approximation (let w be M -dim)

- Taylor series expansion

$$E(w) = E(w') + \nabla E(w')^T (w - w') + \frac{1}{2} (w - w')^T H_E(w') (w - w') + H.O.T.$$

Constant term
Linear (first-order) term
Quadratic (second-order) term

Variable to solve
Current guess
Variation “around” w'

- Rewrite $w = (w' + \Delta w)$

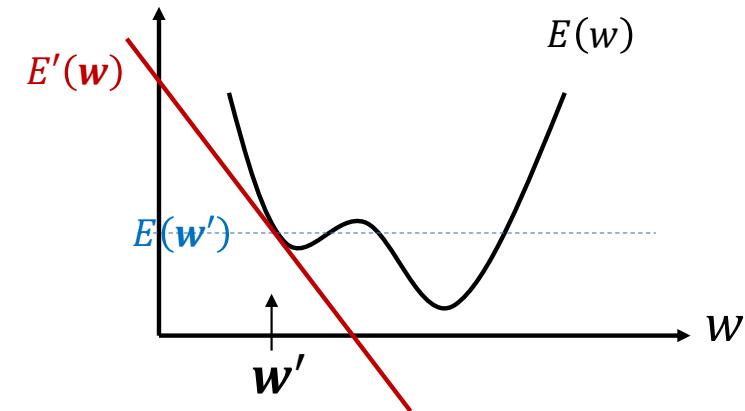
$$E(w' + \Delta w) = E(w') + \nabla E(w')^T \Delta w + \frac{1}{2} \Delta w^T H_E(w') \Delta w + H.O.T.$$

Δw : the movement from w'

Local approximation (let w be M -dim)

- Visualization

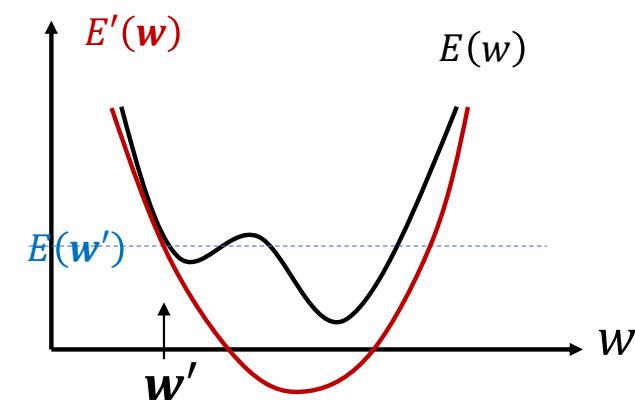
- First order: $E'(w) = E(w') + \nabla E(w')^T (w - w')$



- Second order: $E'(w) = E(w') + \nabla E(w')^T (w - w')$

$$+ \frac{1}{2} (w - w')^T \mathbf{H}_E(w') (w - w')$$

- The approximation $E'(w)$ is quadratic!
 - If $E(w)$ is strictly convex, its Hessian matrix $\mathbf{H}_E(w') > 0$: closed forms to minimize the approximation $E'(w)$



Why does gradient descent work?

- Recall the Taylor expansion

$$E(\mathbf{w}) = E(\mathbf{w}') + \nabla E(\mathbf{w}')^T (\mathbf{w} - \mathbf{w}') + \text{H.O.T.} = E(\mathbf{w}') + \nabla E(\mathbf{w}')^T \Delta\mathbf{w} + \text{H.O.T.}$$

- Approximate $E(\mathbf{w})$ by the linear terms: $E'(\mathbf{w}' + \Delta\mathbf{w}) = E(\mathbf{w}') + \nabla E(\mathbf{w}')^T \Delta\mathbf{w}$
- Which direction, i.e., $\Delta\mathbf{w}$ with $\|\Delta\mathbf{w}\|_2^2 = 1$, can minimize $E'(\mathbf{w}' + \Delta\mathbf{w})$?
 - Answer: direction of $\Delta\mathbf{w} = -\frac{\nabla E(\mathbf{w}')}{\|\nabla E(\mathbf{w}')\|_2}$
 - Interpretation: move in the direction of $-\nabla E(\mathbf{w}')$ is the most efficient for minimization
 - Only works for small $\Delta\mathbf{w}$ (so a scale η) due to the nature of linear Taylor expansion

The GD Algorithm

Initialize parameters:

$$\hat{\mathbf{w}}^{(0)} = [w[1], w[2], \dots, w[M]]^T$$

Error function:

$$E(\mathbf{w})$$

Compute gradients:

$$\nabla E(\hat{\mathbf{w}}^{(t)}) = \begin{bmatrix} \frac{\partial E(\hat{\mathbf{w}}^{(t)})}{\partial w[1]} \\ \vdots \\ \frac{\partial E(\hat{\mathbf{w}}^{(t)})}{\partial w[M]} \end{bmatrix}$$

Update parameters: $\hat{\mathbf{w}}^{(t+1)} = \hat{\mathbf{w}}^{(t)} - \eta^{(t)} \cdot \nabla E(\hat{\mathbf{w}}^{(t)})$

[Repeat until converges]

Learning rate

2-order Taylor expansion approximation

- Taylor series expansion ($\Delta\mathbf{w} = \mathbf{w} - \mathbf{w}'$)

$$E(\mathbf{w}' + \Delta\mathbf{w}) = E(\mathbf{w}') + \nabla E(\mathbf{w}')^T \Delta\mathbf{w} + \frac{1}{2} \Delta\mathbf{w}^T \mathbf{H}_E(\mathbf{w}') \Delta\mathbf{w} + \text{H.O.T.}$$

- Hessian matrix $\mathbf{H}_E(\mathbf{w}')$

$$\mathbf{H}_E(\mathbf{w}') = \begin{bmatrix} \frac{\partial^2 E(\mathbf{w}')}{\partial w[1]\partial w[1]} & \cdots & \frac{\partial^2 E(\mathbf{w}')}{\partial w[1]\partial w[M]} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 E(\mathbf{w}')}{\partial w[M]\partial w[1]} & \cdots & \frac{\partial^2 E(\mathbf{w}')}{\partial w[M]\partial w[M]} \end{bmatrix}$$

- Newton's method

- For any kind of “convex” loss $E(\mathbf{w})$
- Quadratic approximation of $E(\mathbf{w})$ at \mathbf{w}' : ignore H.O.T.
- Minimize the quadratic approximation to get $\Delta\mathbf{w} = -\mathbf{H}_E^{-1}(\mathbf{w}') \nabla E(\mathbf{w}')$
- Update \mathbf{w}' by $\mathbf{w}' + \eta^{(t)} \Delta\mathbf{w}$

- Pros: faster convergence; Cons: $\mathbf{H}_E(\mathbf{w}')$ and $\mathbf{H}_E^{-1}(\mathbf{w}')$ are time consuming

The Newton's method

Initialize parameters:

$$\hat{\mathbf{w}}^{(0)} = [w[1], w[2], \dots, w[M]]^T$$

Error function:

$$E(\mathbf{w})$$

Compute gradients and Hessian:

$$\nabla E(\hat{\mathbf{w}}^{(t)}) = \begin{bmatrix} \frac{\partial E(\hat{\mathbf{w}}^{(t)})}{\partial w[1]} \\ \vdots \\ \frac{\partial E(\hat{\mathbf{w}}^{(t)})}{\partial w[M]} \end{bmatrix}, \quad \mathbf{H}_E(\hat{\mathbf{w}}^{(t)}) = \begin{bmatrix} \frac{\partial^2 E(\hat{\mathbf{w}}^{(t)})}{\partial w[1] \partial w[1]} & \cdots & \frac{\partial^2 E(\hat{\mathbf{w}}^{(t)})}{\partial w[1] \partial w[M]} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 E(\hat{\mathbf{w}}^{(t)})}{\partial w[M] \partial w[1]} & \cdots & \frac{\partial^2 E(\hat{\mathbf{w}}^{(t)})}{\partial w[M] \partial w[M]} \end{bmatrix}$$

Update
parameters:

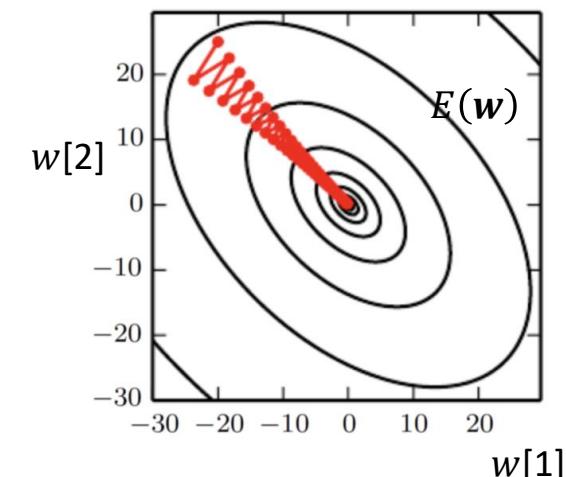
$$\hat{\mathbf{w}}^{(t+1)} = \hat{\mathbf{w}}^{(t)} - \eta^{(t)} \mathbf{H}_E^{-1}(\mathbf{w}') \nabla E(\mathbf{w}')$$

[Repeat until converges]

Learning rate

Summary

- Pros of GD:
 - No expensive matrix operations (i.e., inversion)
 - Non-RSS error functions (if gradient is well formed)
 - Easy to parallelize: $E(\mathbf{w}) = \sum_i \ell(h(\mathbf{x}_i; \mathbf{w}), y_i) \Rightarrow \nabla E(\mathbf{w}) = \sum_i \nabla \ell(h(\mathbf{x}_i; \mathbf{w}), y_i)$
- Cons of GD:
 - Can be misled by local minima and stuck at larger-error valleys
 - Can get stuck on a flat plain/flat region of the error function
 - Can overshoot back and forth (oscillate from side to side)
- Helpers
 - Use multiple random-restarts with random initial parameters
 - Methods exist for selecting appropriate (best) learning rate at each iteration



Summary

- Pros of Newton's method:
 - Non-RSS error functions (if gradient and Hessian are well formed)
 - If \mathbf{w}' is well-chosen, it can converge faster, since 2nd-order Taylor has more information.
- Cons of Newton's method:
 - Hessian and its inverse are time consuming
 - Divergence often happens if the function is flat or almost flat with respect to some dimension. In that case the second derivatives are close to zero, and their inverse becomes very large, resulting in gigantic steps.
- Helpers
 - Compute only the diagonal Hessian: $\Delta\mathbf{w} = -\mathbf{H}_E^{-1}(\mathbf{w}')\nabla E(\mathbf{w}')$
 - GD and then Newton's method

Question

- How many steps do you need to optimize $E(w) = \sum_i (wx_i - y_i)^2$?
 - GD?
 - Newton's method?

Today

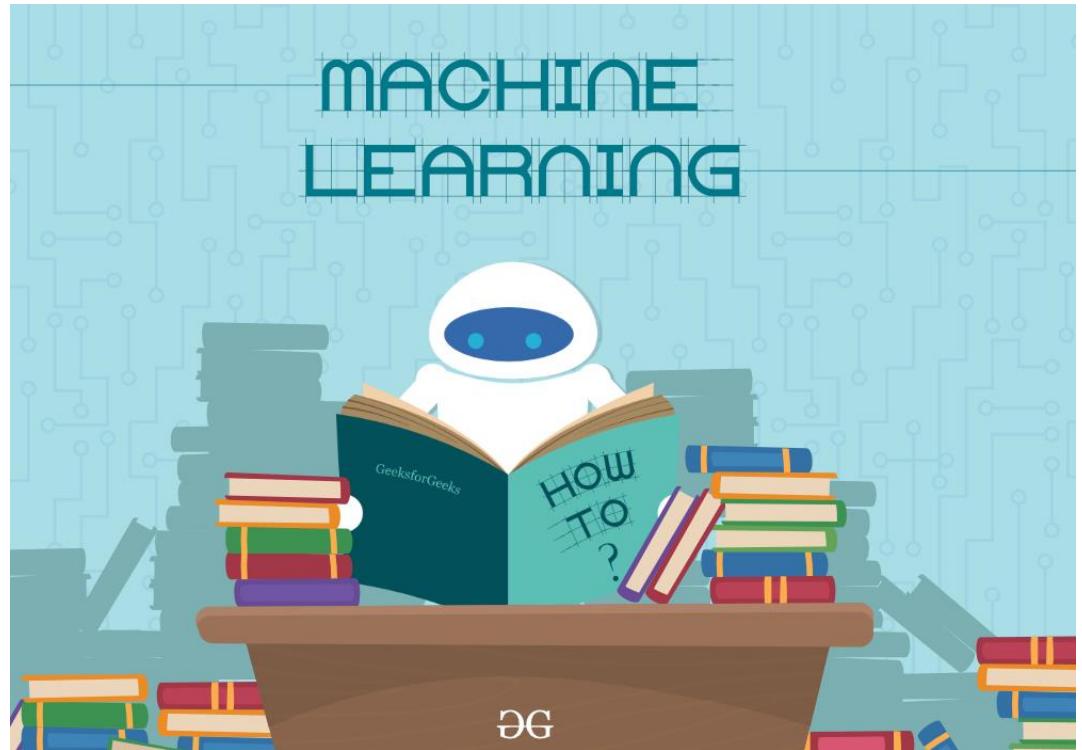
Optimization

- Gradient descent (continued)
- Newton's method

Linear classifier

- Definition
- First algorithm: perceptron
- Reading: pocket algorithm

Probability refresher (self reading)



“Binary Linear” classification

- Predicting a binary class label
 - Yes/No, 1/0, +1/-1
- Setup
 - $D_{tr} = \{(x_n \in \mathbb{R}^D, y_n \in \{-1, +1\})\}_{n=1}^N$
- Models (hypothesis class = hyperplanes):
 - $\hat{y} = h(x) = \text{sign}(\mathbf{w}^T \mathbf{x} + b) = \text{sign}(\mathbf{w}^T \mathbf{x} + w[0]) = \text{sign}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})$
 - $\tilde{\mathbf{w}} = [\mathbf{w}[0], w[1], \dots, w[D]]^T \in \mathbb{R}^{D+1}; \tilde{\mathbf{x}} = [1, x[1], \dots, x[D]]^T \in \mathbb{R}^{D+1}$
 - \mathbf{w} (or $\tilde{\mathbf{w}}$): weight vector; b : bias
 - Assumption: Linear separation between input features of different classes

Linear classification

- Data visualization and task illustration

- What if $b = 0$?

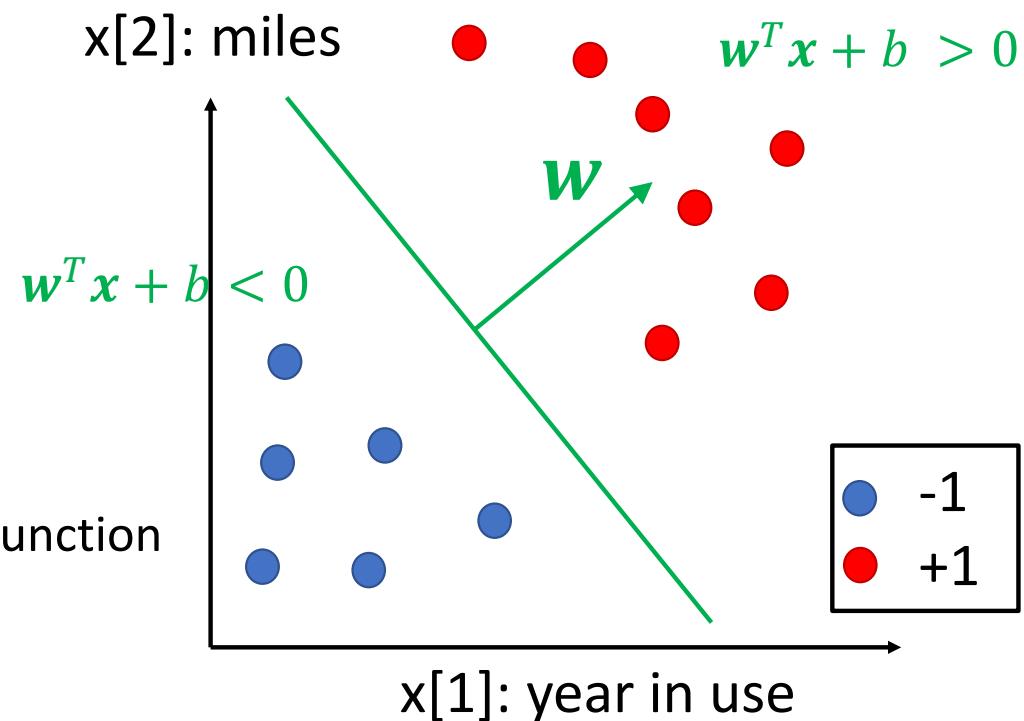
- Meaning:

- If $w[1] * \text{year} + w[2] * \text{miles} > -b$, do ...

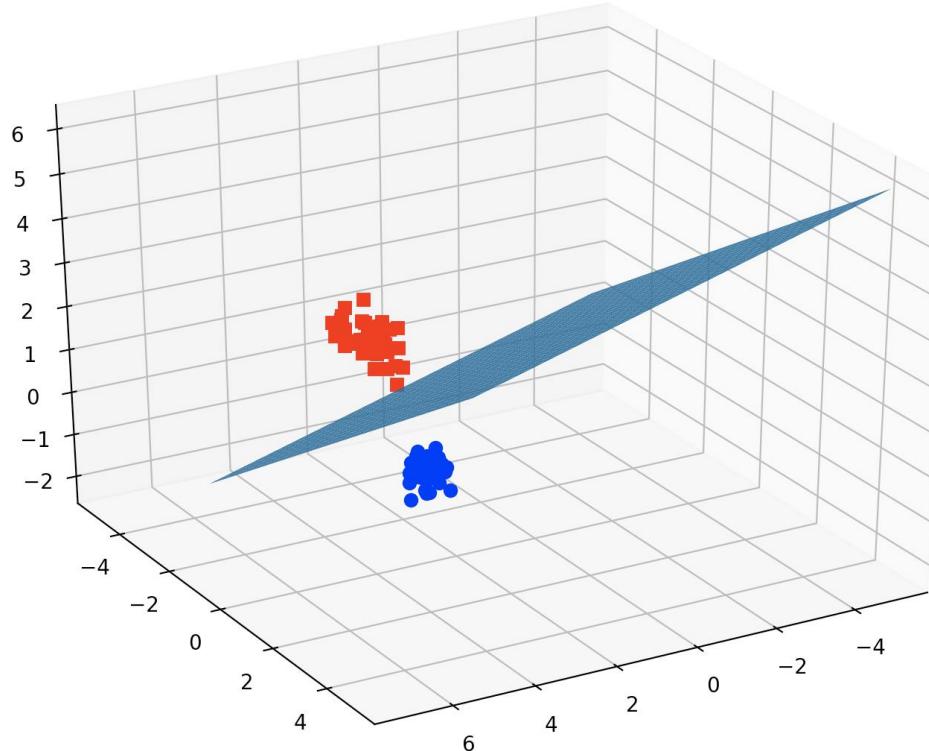
- How to measure the loss?

- 0-1 loss: $\ell(y, y') = 1[y \neq y']$

- $1[\text{True}] = 1, 1[\text{False}] = 0$ is the indicator function



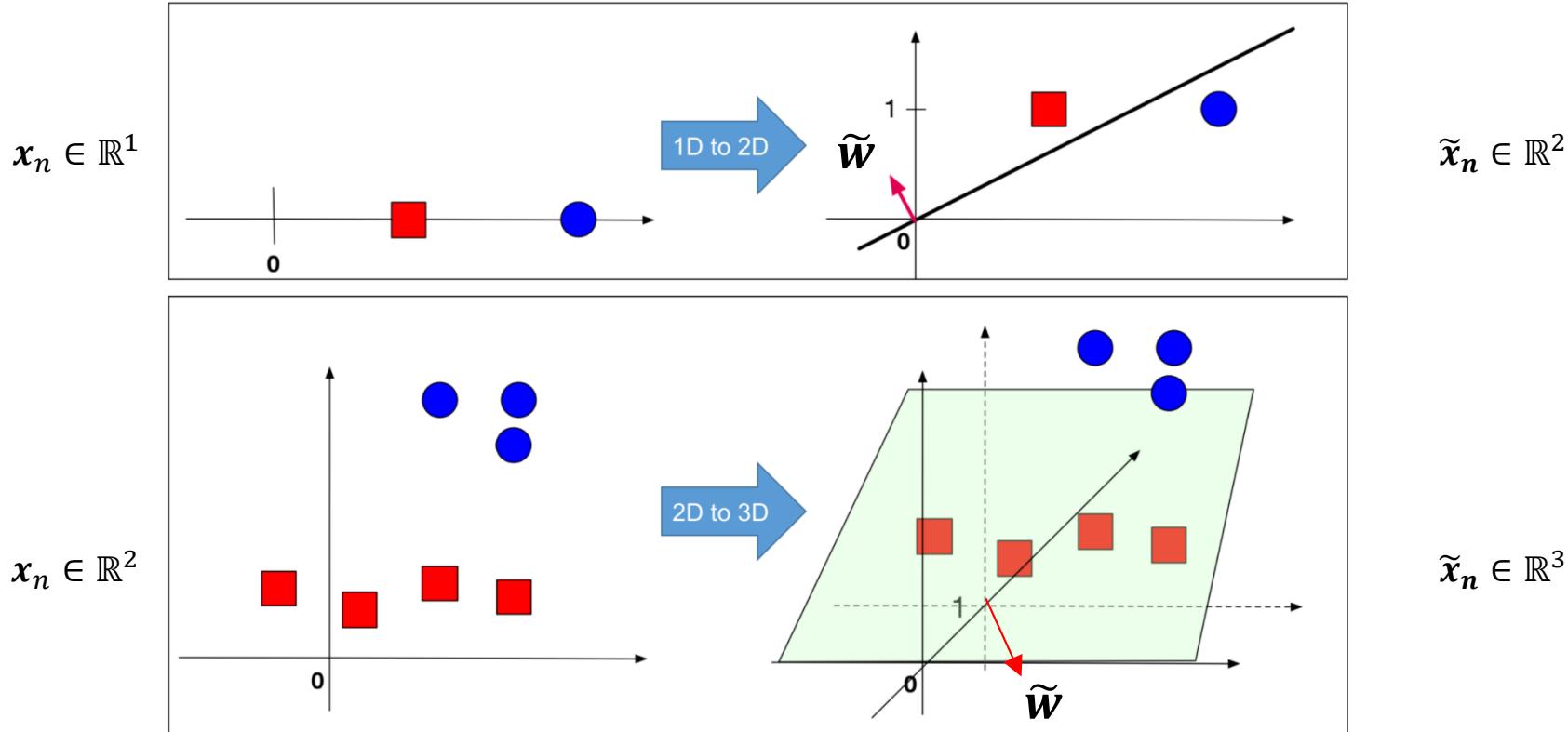
Linear classification: beyond 2-D input



Learning goal

- Ideally, minimizing the generalized classification error as much as possible
 - But we only have training data
 - So, let's try empirical risk minimization
- **0-1 classification error**
 - $E(\tilde{\mathbf{w}}) = \sum_i \mathbf{1}[y_i \neq h(\mathbf{x}_i; \tilde{\mathbf{w}})] = \sum_i \mathbf{1}[y_i \neq \text{sign}(\mathbf{w}^T \mathbf{x}_i + b)] = \sum_i \mathbf{1}[y_i \neq \text{sign}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i)]$
 - $E(\tilde{\mathbf{w}}) = \sum_i \mathbf{1}[y_i (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i) \leq 0]$
 - Is this easy to minimize?
- Others:
 - Surrogate losses as “differentiable” upper bounds for optimization

Visualization of the concatenation



Terminology

- **Classification accuracy**

- $\frac{1}{N} \sum_i \mathbf{1}[y_i = h(\mathbf{x}_i; \tilde{\mathbf{w}})]$
 - If it is on training/test data, then called training/test accuracy.

- **Note:**

- Later when we introduce other ML algorithms, we may not use 0-1 loss!
 - Nevertheless, when we calculate the accuracy, we use the formula above.

Today

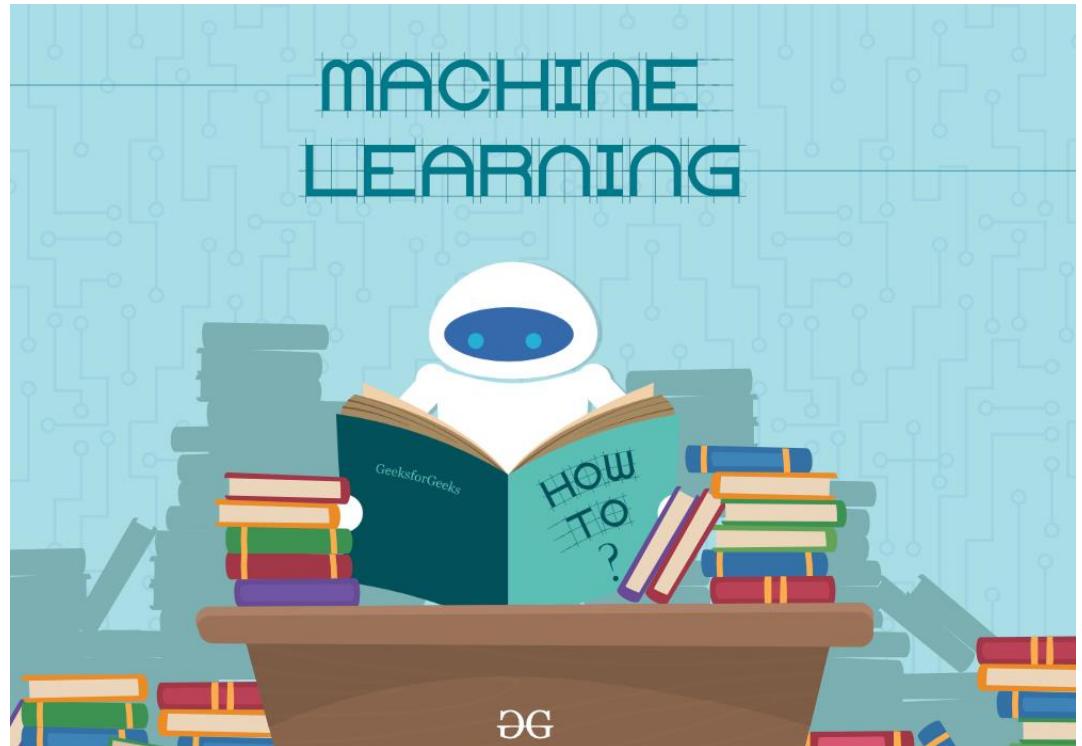
Optimization

- Gradient descent (continued)
- Newton's method

Linear classifier

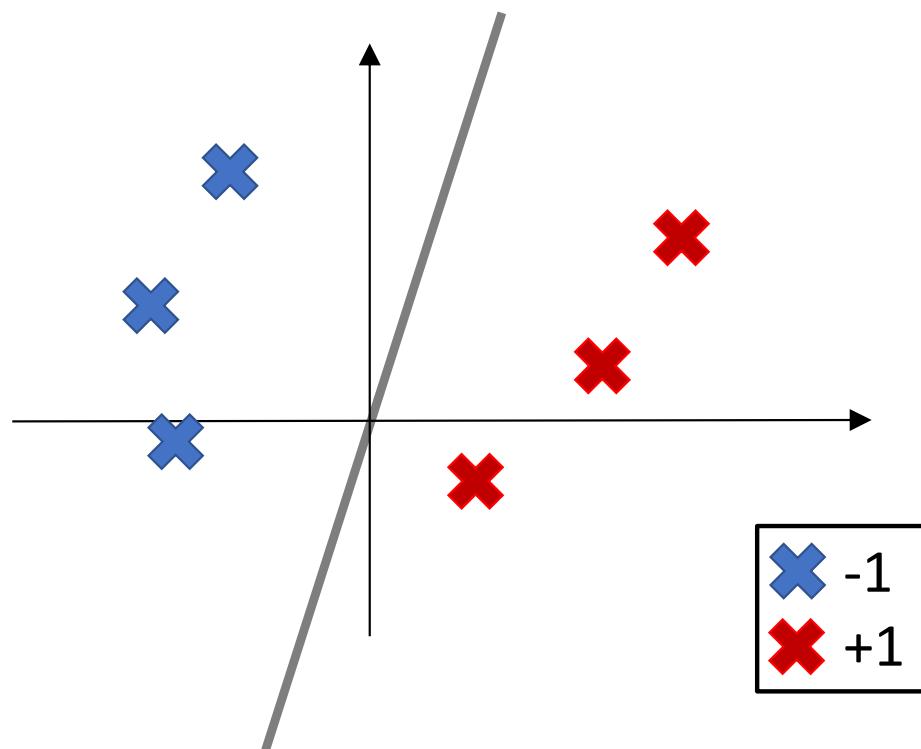
- Definition
- First algorithm: perceptron
- Reading: pocket algorithm

Probability refresher (self reading)



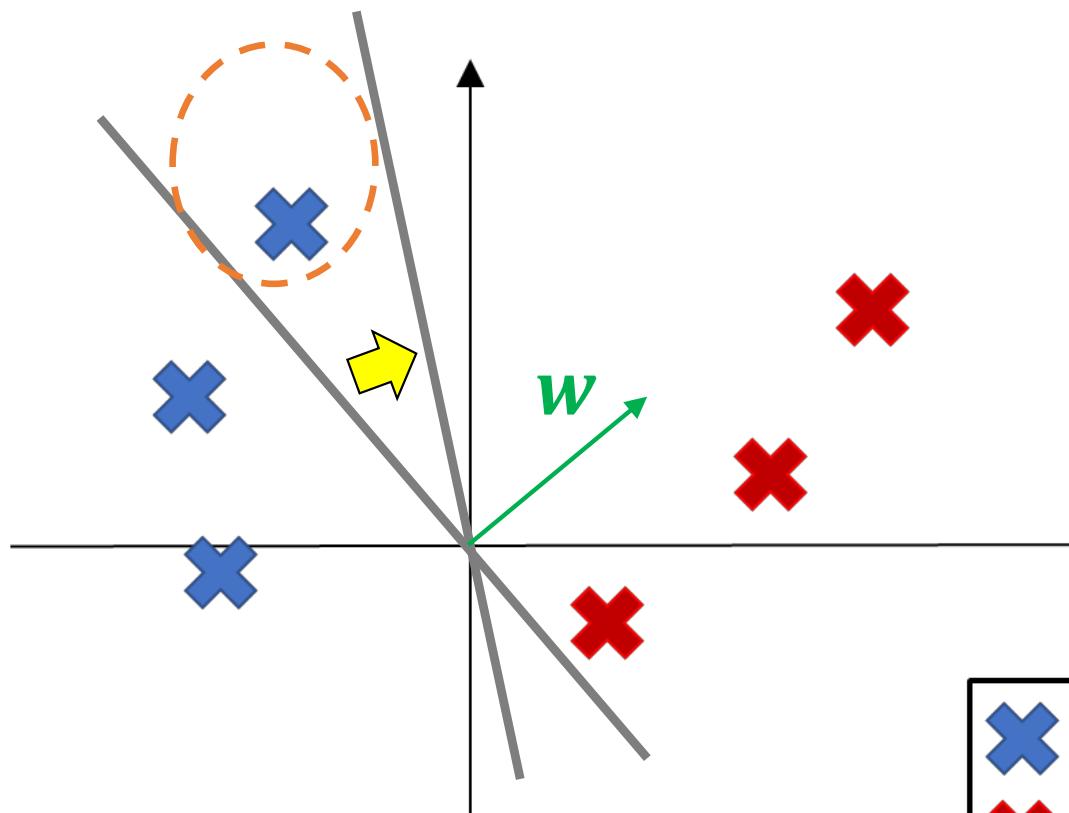
Perceptron Algorithm

- An algorithm to learn a binary classification by minimizing the 0-1 loss



Perceptron Algorithm

- Concept:



| |
|---------------------|
| $\text{X} \quad -1$ |
| $\text{X} \quad +1$ |

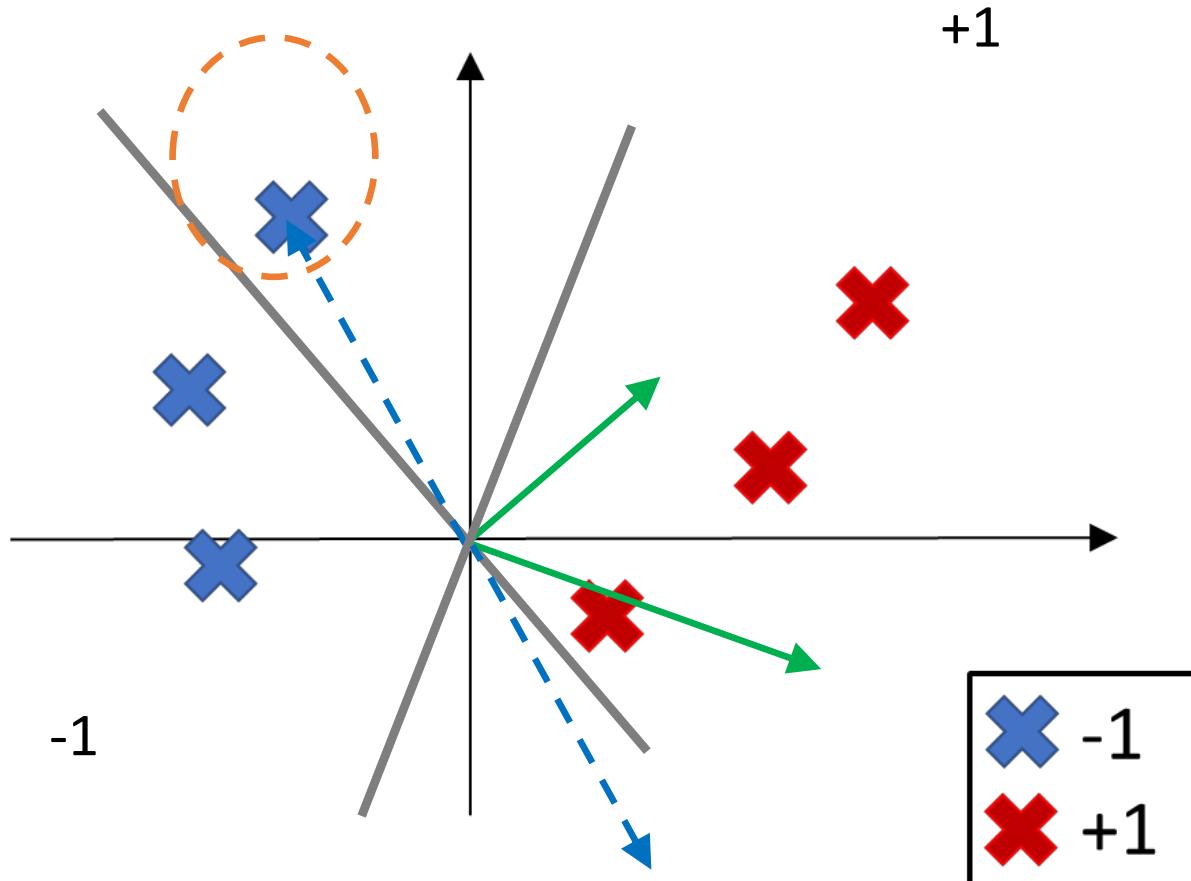
Perceptron Algorithm

- Let $x_n \in \mathbb{R}^D$ and $w \in \mathbb{R}^D$, $y \in \{-1,1\}$
 - Assume that a dummy variable 1 is included in x_n , and b is included in w
- Initialize weight vector $w = \mathbf{0}$
- Loop for T iterations (or till converge; i.e., no misclassification)
 - Loop for all training examples x_n (random order!)
 - Predict $\hat{y}_n = \text{sign}(w^T x_n)$
 - If $\hat{y}_n \neq y_n$
 - Update: $w \leftarrow w + \eta(y_n x_n)$
- Note: η is the learning rate (step size) and $\eta \geq 0$

Perceptron Algorithm

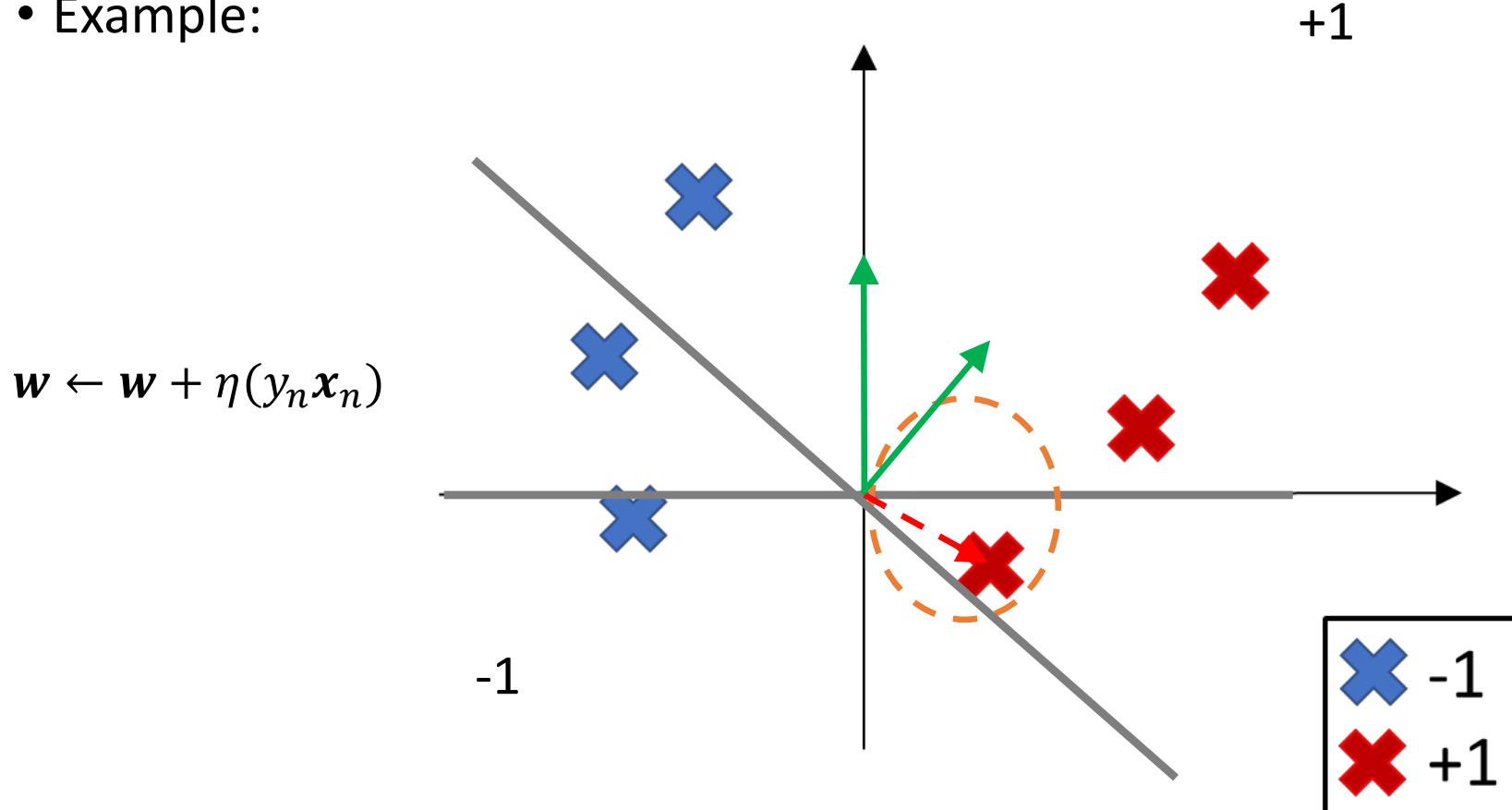
- Example:

$$w \leftarrow w + \eta(y_n x_n)$$



Perceptron Algorithm

- Example:



Perceptron Algorithm

- Why does it work?
 - If $\hat{y}_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n) = -1$ but $y_n = +1$
 - $(\mathbf{w} + \eta(y_n \mathbf{x}_n))^T \mathbf{x}_n = \mathbf{w}^T \mathbf{x}_n + \eta \textcolor{green}{y_n} \textcolor{blue}{x_n}^T \mathbf{x}_n \geq \mathbf{w}^T \mathbf{x}_n$
 - If $\hat{y}_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n) = +1$ but $y_n = -1$
 - $(\mathbf{w} + \eta(y_n \mathbf{x}_n))^T \mathbf{x}_n = \mathbf{w}^T \mathbf{x}_n + \eta \textcolor{green}{y_n} \textcolor{blue}{x_n}^T \mathbf{x}_n \leq \mathbf{w}^T \mathbf{x}_n$

Questions?



THE OHIO STATE UNIVERSITY

Perceptron Algorithm: online/streaming version

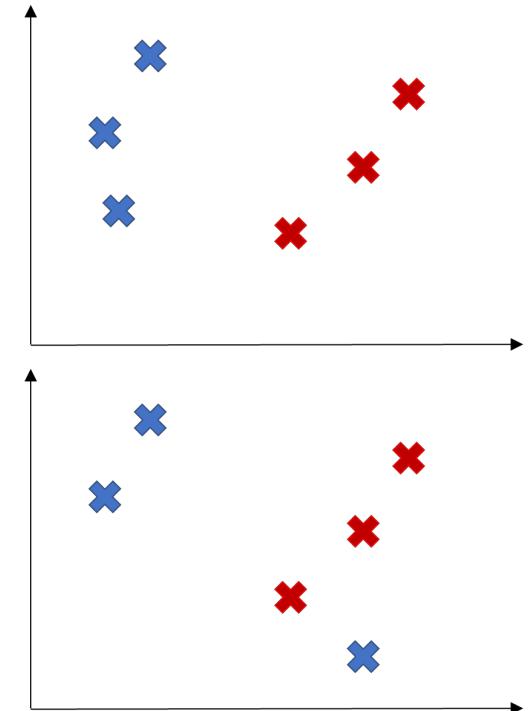
- Let $x \in \mathbb{R}^D$ and $w \in \mathbb{R}^D$, $y \in \{-1, 1\}$
 - Assume that a dummy variable 1 is included in x_n , and b is included in w
 - Initialize weight vector $w = \mathbf{0}$
 - **Online version:** Get one training data (x, y)
 - Predict $\hat{y} = \text{sign}(w^T x)$
 - If $\hat{y} \neq y$
 - **Update:** $w \leftarrow w + \eta(yx)$

Perceptron Algorithm

- Convergence (in terms of the training error):
 - If the data is linear separable

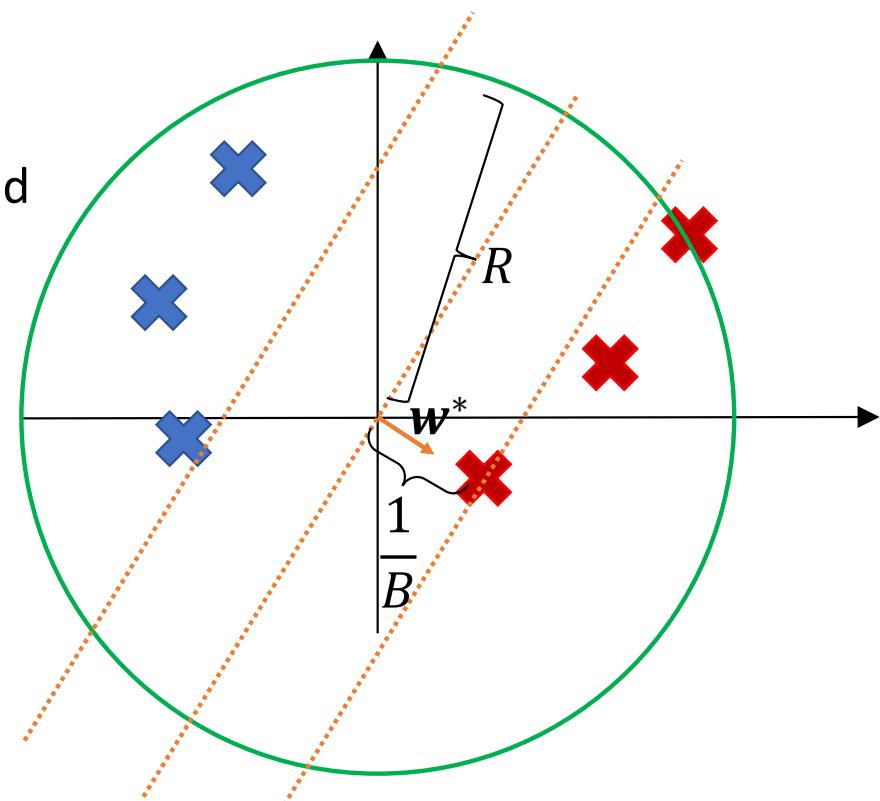
- Not convergence:
 - If the data is not linear separable
 - Not good ... any idea to resolve it?

- **The pocket algorithm:** keep the best solution so far “in the pocket”
 - “Best” in terms of the training error



Proof of convergence

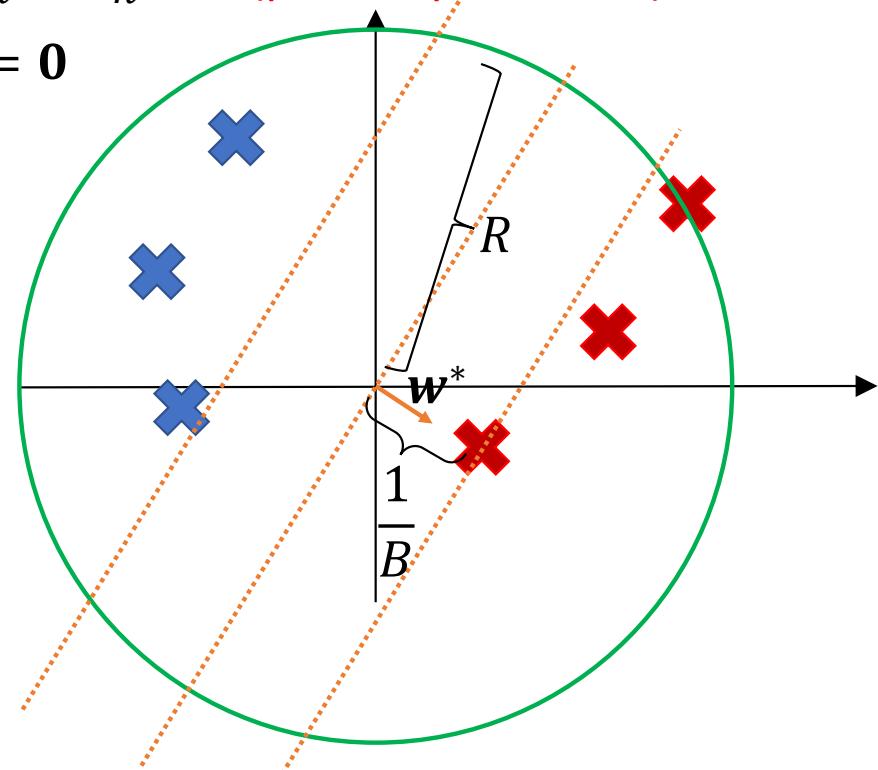
- If the training set $D_{tr} = \{(x_n \in \mathbb{R}^D, y_n \in \{-1, +1\})\}_{n=1}^N$ is linearly separable
 - $R = \max_n \|x_n\|$
 - $B = \min_n \|\mathbf{w}\|, s.t. \forall n, y_n \mathbf{w}^T \mathbf{x}_n \geq 1$
 - Stop after performing at most $(RB)^2$ updated
- Intuition:
 - The more separable the data are
 - The faster the convergence



Proof of convergence

- Proof sketch

- Let \mathbf{w}^* be a solution of $\min_n \|\mathbf{w}\|_2$, s.t. $\forall n, y_n \mathbf{w}^T \mathbf{x}_n \geq 1$ (perfect prediction)
- Let $\mathbf{w}^{(t)}$ be the \mathbf{w} after t updates, and $\mathbf{w}^{(0)} = \mathbf{0}$
- Let $\eta = 1$
- Show that $\frac{\mathbf{w}^{(t)T} \mathbf{w}^*}{\|\mathbf{w}^{(t)}\|_2 \|\mathbf{w}^*\|_2} \geq \frac{\sqrt{t}}{RB}$
 - The left-hand side cannot go over 1
 - $t \leq (RB)^2$
- Intuition:
 - every update moves $\mathbf{w}^{(t)}$ closer to \mathbf{w}^*
- Details:
 - See UML: 9-9.1.2



Today

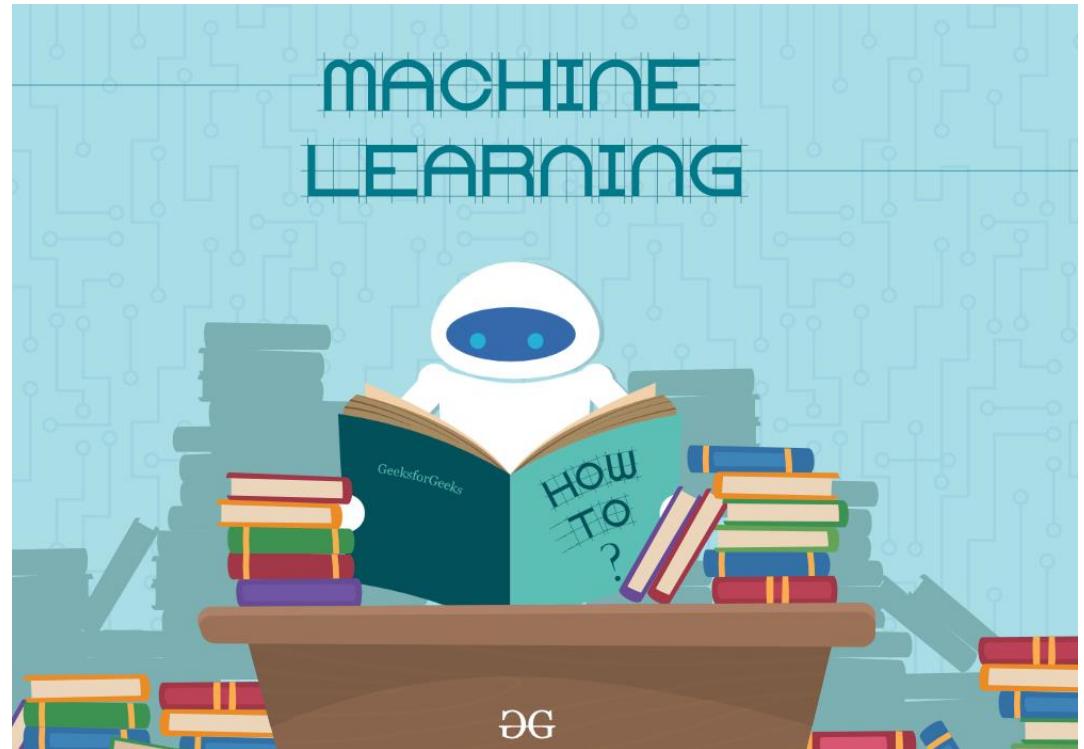
Optimization

- Gradient descent (continued)
- Newton's method

Linear classifier

- Definition
- First algorithm: perceptron
- Reading: pocket algorithm

Probability refresher (self reading)



Reading: Pocket algorithm

- **Training data:** $D_{tr} = \{ (x_i \in \mathbb{R}^D, y_i \in \{+1, -1\}) \}_{i=1}^N$
- **Model:** $\text{sign}(wx + b) = \text{sign}(\tilde{w}^T \tilde{x})$
- **Goal:** Make Perceptron stable when the data is not linearly separable
- **How?**
 - After every epoch (every iteration through all the training data), you compare the current solution \tilde{w} to \tilde{w}^{best} in terms of their accuracy on the training data
 - If \tilde{w} is better than \tilde{w}^{best} , update \tilde{w}^{best} by \tilde{w} . That is, \tilde{w}^{best} keeps the best solution in terms of the training accuracy so far.

Reading: Pocket algorithm

- Initialize $\tilde{\mathbf{w}}$ and $\tilde{\mathbf{w}}^{\text{best}} = \mathbf{0}$
- For $t = 1: T$
 - Loop for all training examples \tilde{x}_n (random order!)
 - Predict $\hat{y}_n = \text{sign}(\tilde{\mathbf{w}}^T \tilde{x}_n)$
 - If $\hat{y}_n \neq y_n$
 - Update: $\tilde{\mathbf{w}} \leftarrow \tilde{\mathbf{w}} + \eta(y_n \tilde{x}_n)$
 - Evaluate $\tilde{\mathbf{w}}$ on the “training data” and calculate the training accuracy
 - If training accuracy by $\tilde{\mathbf{w}}$ is “higher” than the training accuracy by $\tilde{\mathbf{w}}^{\text{best}}$
 - $\tilde{\mathbf{w}}^{\text{best}} \leftarrow \tilde{\mathbf{w}}$
- Output $\tilde{\mathbf{w}}^{\text{best}}$

This is an epoch: seeing all the training data once

Reading: Pocket algorithm

- Initialize $\tilde{\mathbf{w}}$ and $\tilde{\mathbf{w}}^{\text{best}} = \mathbf{0}$
- For $t = 1:T$
 - Loop for all training examples $\tilde{\mathbf{x}}_n$ (random order!)
 - Predict $\hat{y}_n = \text{sign}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)$
 - If $\hat{y}_n \neq y_n$
 - Update: $\tilde{\mathbf{w}} \leftarrow \tilde{\mathbf{w}} + \eta(y_n \tilde{\mathbf{x}}_n)$
 - Evaluate $\tilde{\mathbf{w}}$ on the “training data” and calculate the training accuracy
 - If training accuracy by $\tilde{\mathbf{w}}$ is “higher” than the training accuracy by $\tilde{\mathbf{w}}^{\text{best}}$
 - $\tilde{\mathbf{w}}^{\text{best}} \leftarrow \tilde{\mathbf{w}}$
- Output $\tilde{\mathbf{w}}^{\text{best}}$

Here you compare $\tilde{\mathbf{w}}$ to $\tilde{\mathbf{w}}^{\text{best}}$ and make an update on $\tilde{\mathbf{w}}^{\text{best}}$ if needed

Reading: Pocket algorithm

- Initialize $\tilde{\mathbf{w}}$ and $\tilde{\mathbf{w}}^{\text{best}} = \mathbf{0}$
- For $t = 1:T$
 - Loop for all training examples $\tilde{\mathbf{x}}_n$ (random order!)
 - Predict $\hat{y}_n = \text{sign}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)$
 - If $\hat{y}_n \neq y_n$
 - Update: $\tilde{\mathbf{w}} \leftarrow \tilde{\mathbf{w}} + \eta(y_n \tilde{\mathbf{x}}_n)$
 - Evaluate $\tilde{\mathbf{w}}$ on the “training data” and calculate the training accuracy
 - If training accuracy by $\tilde{\mathbf{w}}$ is “higher” than the training accuracy by $\tilde{\mathbf{w}}^{\text{best}}$
 - $\tilde{\mathbf{w}}^{\text{best}} \leftarrow \tilde{\mathbf{w}}$
- Output $\tilde{\mathbf{w}}^{\text{best}}$

Note that, this step does not update $\tilde{\mathbf{w}}$

Reading: Pocket algorithm

- Initialize $\tilde{\mathbf{w}}$ and $\tilde{\mathbf{w}}^{\text{best}} = \mathbf{0}$
- For $t = 1:T$
 - Loop for all training examples $\tilde{\mathbf{x}}_n$ (random order!)
 - Predict $\hat{y}_n = \text{sign}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)$
 - If $\hat{y}_n \neq y_n$
 - Update: $\tilde{\mathbf{w}} \leftarrow \tilde{\mathbf{w}} + \eta(y_n \tilde{\mathbf{x}}_n)$
 - Evaluate $\tilde{\mathbf{w}}$ on the “training data” and calculate the training accuracy
 - If training accuracy by $\tilde{\mathbf{w}}$ is “higher” than the training accuracy by $\tilde{\mathbf{w}}^{\text{best}}$
 - $\tilde{\mathbf{w}}^{\text{best}} \leftarrow \tilde{\mathbf{w}}$
- Output $\tilde{\mathbf{w}}^{\text{best}}$

Finally, you output $\tilde{\mathbf{w}}^{\text{best}}$

Today

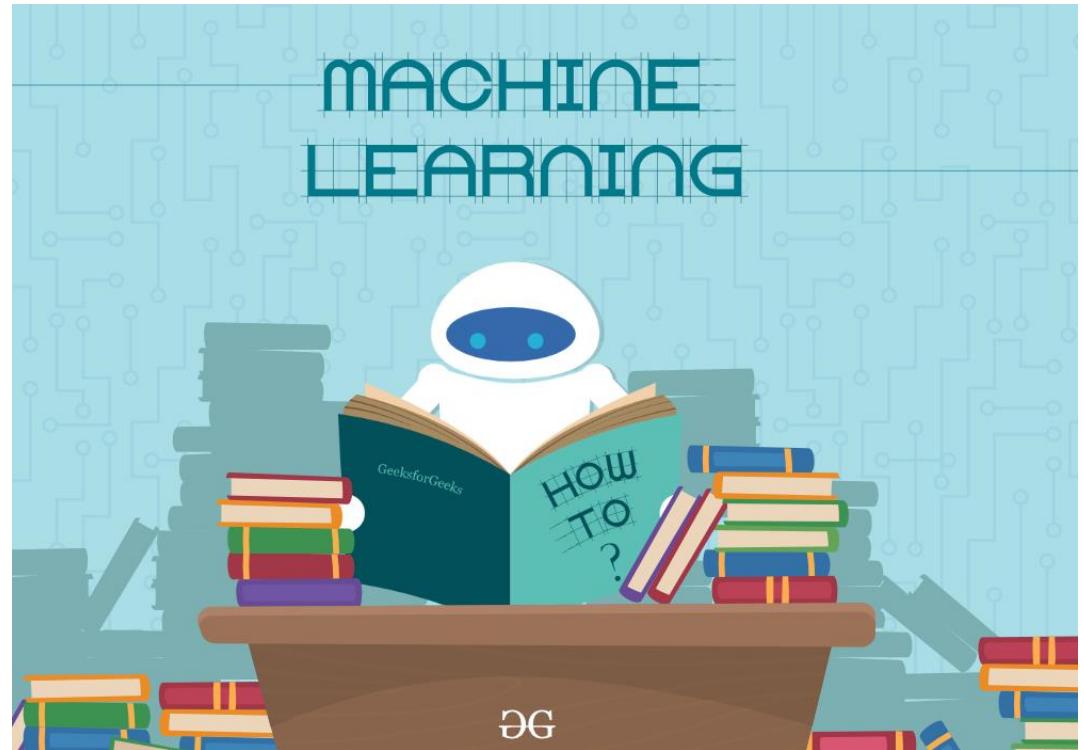
Optimization

- Gradient descent (continued)
- Newton's method

Linear classifier

- Definition
- First algorithm: perceptron
- Reading: pocket algorithm

Probability refresher (self reading)

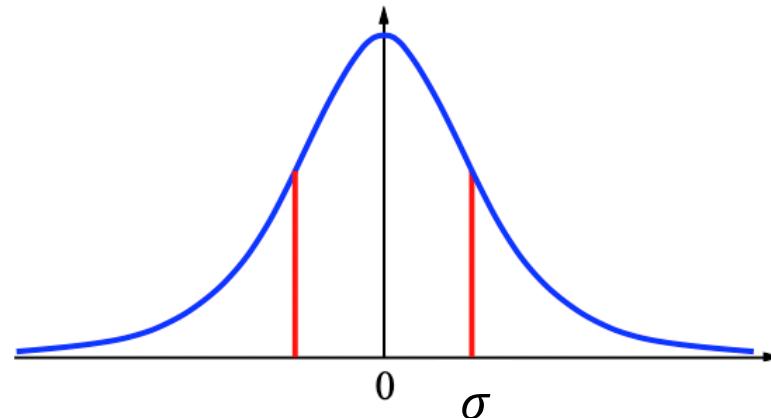


Probability (very high-level)

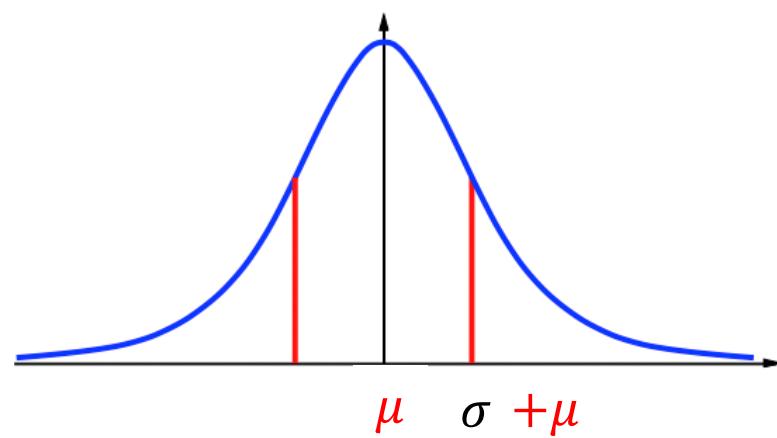
- Let X be a random variable and x be a particular outcome value
- If the sample space \mathcal{X} (where x can take values from) is discrete
 - $P(X = x)$ is the probability mass function: the probability of $X = x$
 - $\sum_{x \in \mathcal{X}} P(X = x) = 1, 0 \leq P(X = x) \leq 1$
- If the sample space \mathcal{X} (where x can take values from) is continuous
 - $p(X = x)$ is the probability density function: $P(b \geq X \geq a) = \int_a^b p(X = x) dx$
 - $\int p(X = x) dx = 1, p(X = x) \geq 0$
- When there is no confusion, X will be ignored; P and p are interchangeable

Gaussian density

- $p(X = \textcolor{blue}{x}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\textcolor{blue}{x}^2/2\sigma^2}$



- $p(X = \textcolor{blue}{x}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(\textcolor{blue}{x}-\textcolor{red}{\mu})^2/2\sigma^2}$



Probability (very high-level)

- Expectation/mean

- $\circ \mu = E_P[X] = \sum_{x \in \mathcal{X}} P(X = x) \textcolor{red}{x}; \quad E_P[X] = E_{x \sim P}[x]$
- $\circ \mu = E_p[X] = \int p(X = x) \textcolor{red}{x} dx; \quad E_p[X] = E_{x \sim p}[x]$

- Variance

- $\circ E_P[(X - \mu)^2] = \sum_{x \in \mathcal{X}} P(X = x) (\textcolor{red}{x} - \mu)^2$
- $\circ E_p[(X - \mu)^2] = \int p(X = x) (\textcolor{red}{x} - \mu)^2 dx$

Probability (very high-level)

- Joint probability (the starting point):
 - $P(X = x, Y = y); \quad \sum_{x \in \mathcal{X}, y \in \mathcal{Y}} P(X = x, Y = y) = 1$
 - $p(X = x, Y = y); \quad \int p(X = x, Y = y) dx dy = 1$
- Marginal probability:
 - $P(X = x) = \sum_{y \in \mathcal{Y}} P(X = x, Y = y)$
 - $p(X = x) = \int p(X = x, Y = y) dy$
- Conditional probability:
 - $P(X = x | Y = y) = \frac{P(X=x, Y=y)}{P(Y=y)} = \frac{P(X=x, Y=y)}{\sum_{x \in \mathcal{X}} P(X=x, Y=y)}$
 - $p(X = x | Y = y) = \frac{p(X=x, Y=y)}{p(Y=y)} = \frac{p(X=x, Y=y)}{\int p(X=x, Y=y) dx}$

Probability (very high-level)

- Chain rule:
 - $P(X = x, Y = y) = P(X = x)P(Y = y|X = x) = P(Y = y)P(X = x|Y = y)$
 - $P(X_1 = x_1, \dots, X_N = x_N) = P(x_1)P(x_2|x_1) \dots P(x_N|x_1, \dots, x_{N-1}) = \prod_i P(x_i|x_1, \dots, x_{i-1})$
- Bayes' rule
 - $P(X = x|Y = y) = \frac{P(X=x, Y=y)}{P(Y=y)} = \frac{P(X=x)P(Y=y|X=x)}{P(Y=y)}$
 - $P(Y = y|X = x) = \frac{P(X=x, Y=y)}{P(X=x)} = \frac{P(Y=y)P(X=x|Y=y)}{P(X=x)}$
 - Usage: reverse the cause and effect

Probability (very high-level)

- Independence (not always)
 - $P(X = x, Y = y) = P(X = x)P(Y = y)$
- Conditional independence (not always)
 - $P(X = x, Y = y|Z = z) = P(X = x|Z = z)P(Y = y|Z = z)$
 - Imply $P(X = x|Z = z, Y = y) = P(X = x|Z = z)$
 - Imply $P(Y = y|Z = z, X = x) = P(Y = y|Z = z)$
- In machine learning, inferring if random variables are independent, from data, is very challenging
- Most of the time, we “assume” variables are (conditional) independent to simplify the model and reduce the number of parameters to learn

Some terminology

- **Probabilistic inference:**
 - Derive the (conditional or marginal) probability of one or more random variables taking a specific value or set of values
 - **Example:**
 - Suppose $P(X = x, Y = y)$ is known
 - Derive $P(X = x|Y = \textcolor{red}{y})$, $P(X = x)$, $\operatorname{argmax}_x P(X = x|Y = y)$
- **Learning/estimation:**
 - How to know the parameters within $P(X = x, Y = y)$?
 - How to know what probabilistic models to use for $P(X = x, Y = y)$?

What are probability models?

- Mathematical representation of probabilities
 - Bernoulli distribution (2 discrete assignments)
 - Categorical distribution (multiple discrete assignments)
 - Gaussian distribution (continuous assignments)
- Probability models “often” have parameters
 - To set by domain experts
 - **To be estimated from data: MLE, MAP, and Bayesian estimates**

The probability of win
is 40%; loss is 60%



CSE 5523: Estimating Probability from Data



THE OHIO STATE UNIVERSITY

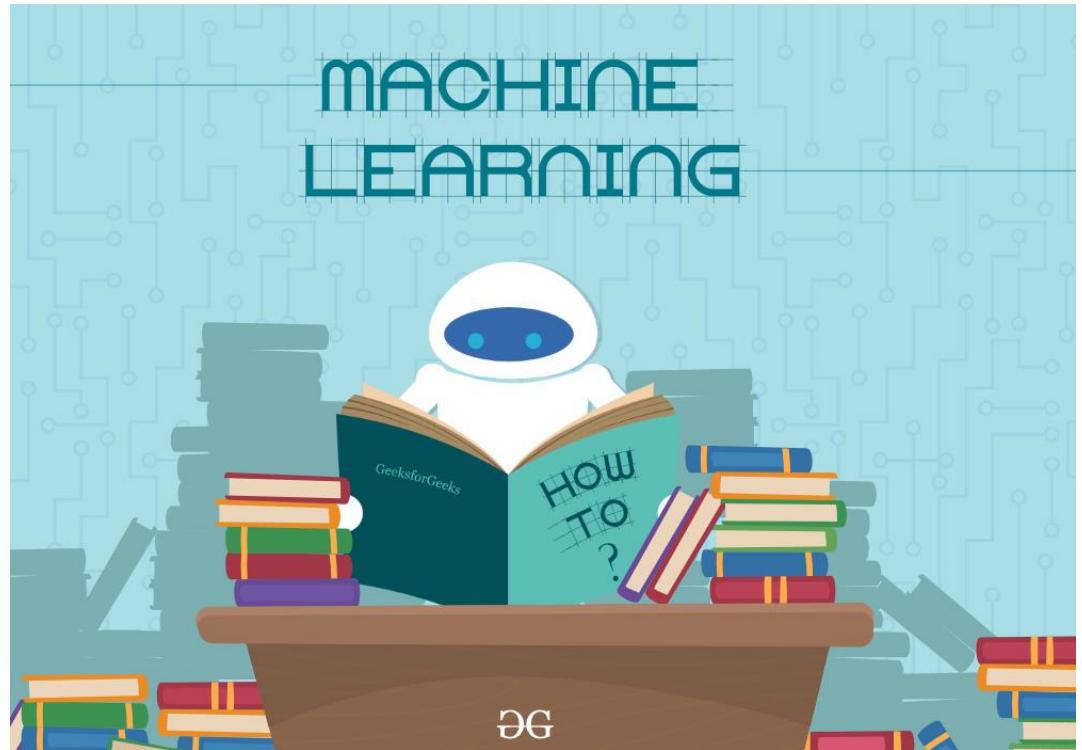
HW

- HW-1 is due next Tuesday at midnight (23:59 ET).
- HW-2 will be released next Tuesday by midnight, and due in two weeks (10/1).

Today – Part 1

Linear classifier

- Review
- Perceptron (continued)

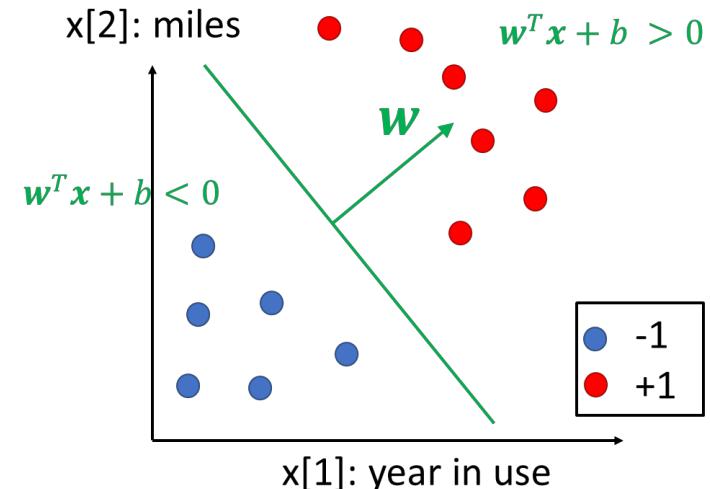


“Binary Linear” classification

- Predicting a binary class label
 - Yes/No, 1/0, +1/-1
- Setup
 - $D_{tr} = \{(x_n \in \mathbb{R}^D, y_n \in \{-1, +1\})\}_{n=1}^N$

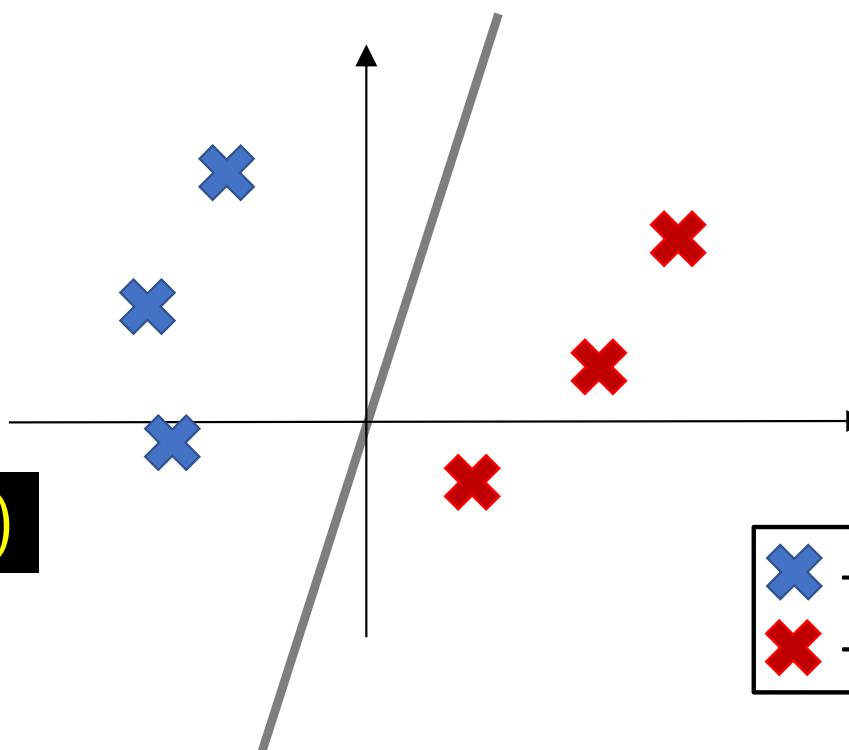
- Models (hypothesis class = hyperplanes):

- $\hat{y} = h(x) = \text{sign}(\mathbf{w}^T \mathbf{x} + b) = \text{sign}(\mathbf{w}^T \mathbf{x} + w[0]) = \text{sign}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})$
- $\tilde{\mathbf{w}} = [w[0], w[1], \dots, w[D]]^T \in \mathbb{R}^{D+1}; \tilde{\mathbf{x}} = [1, x[1], \dots, x[D]]^T \in \mathbb{R}^{D+1}$



Perceptron Algorithm

- An algorithm to learn a binary classification by minimizing the 0-1 loss



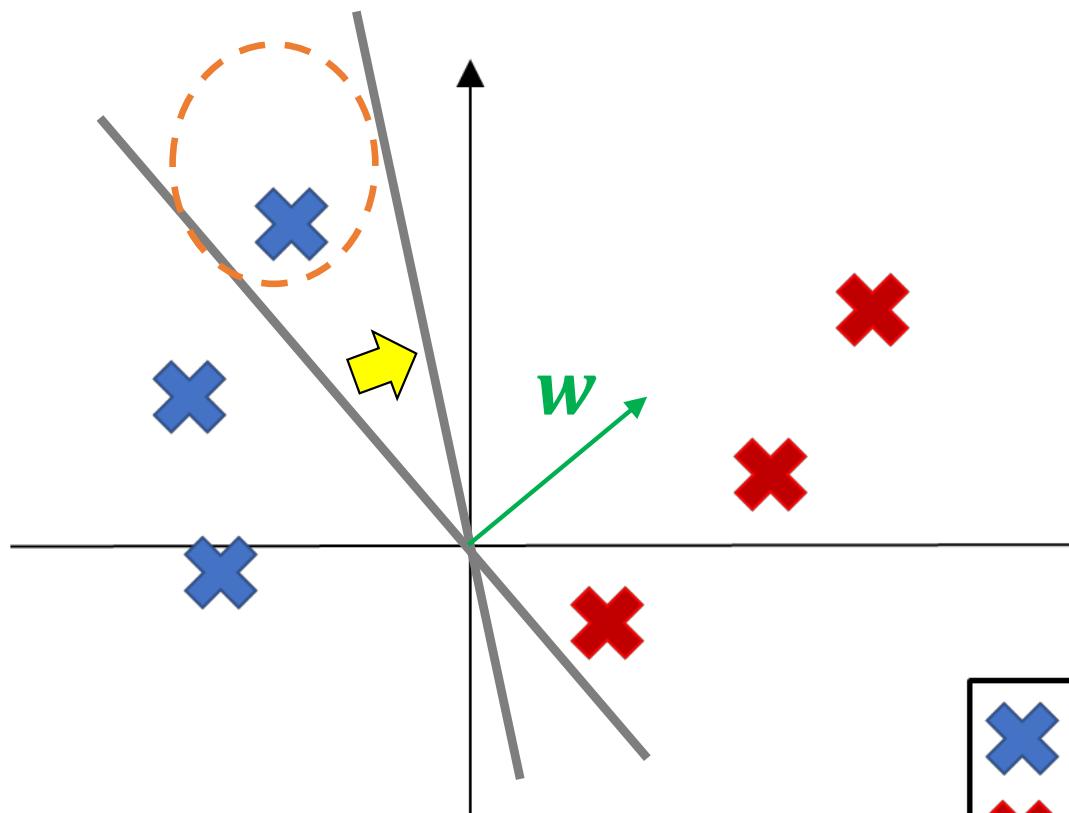
$$\text{sign}(\mathbf{w}^T \mathbf{x})$$

0-1 classification error

$$\begin{aligned} E(\tilde{\mathbf{w}}) &= \sum_i \mathbf{1}[y_i \neq h(\mathbf{x}_i; \tilde{\mathbf{w}})] \\ &= \sum_i \mathbf{1}[y_i \neq \text{sign}(\mathbf{w}^T \mathbf{x}_i + b)] \\ &= \sum_i \mathbf{1}[y_i \neq \text{sign}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i)] \\ E(\tilde{\mathbf{w}}) &= \sum_i \mathbf{1}[y_i (\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i) \leq 0] \end{aligned}$$

Perceptron Algorithm

- Concept:



| |
|---------------------|
| $\text{X} \quad -1$ |
| $\text{X} \quad +1$ |

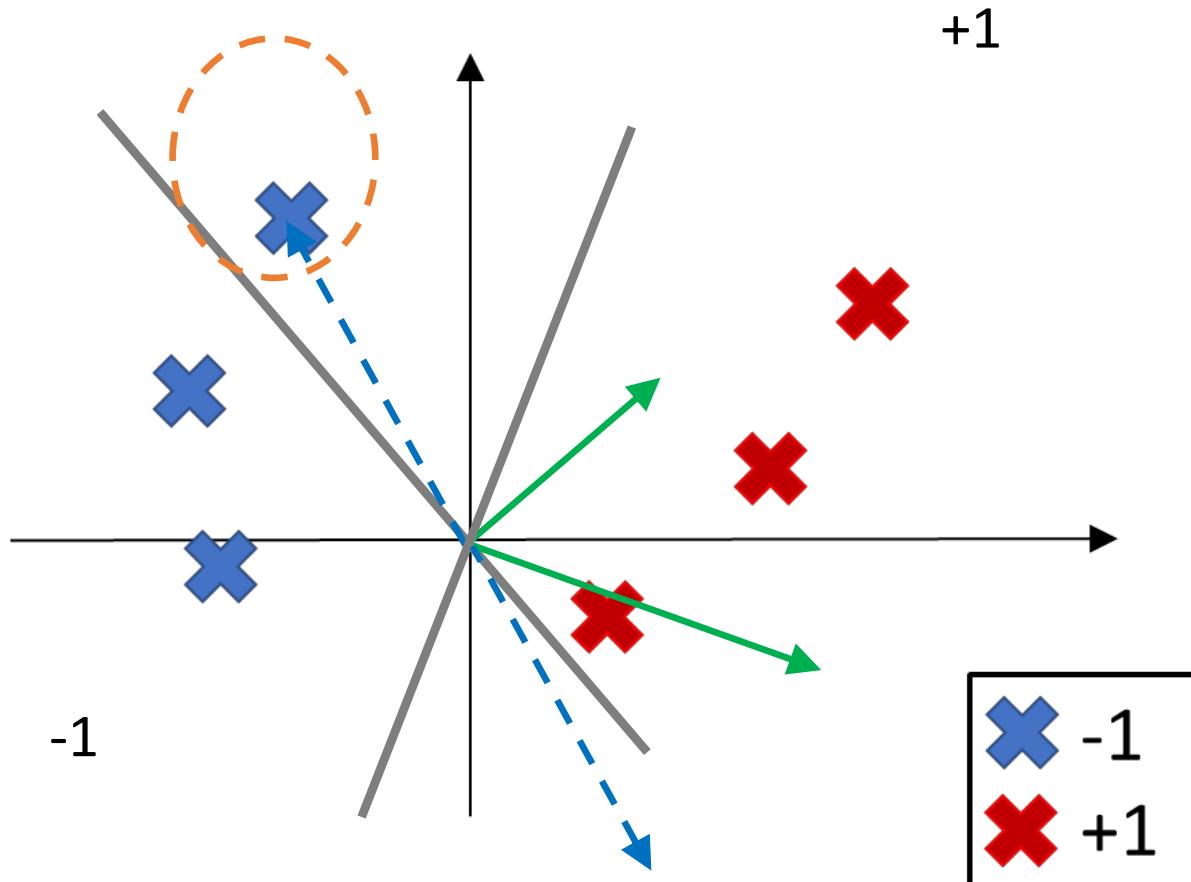
Perceptron Algorithm

- Let $x_n \in \mathbb{R}^D$ and $w \in \mathbb{R}^D$, $y \in \{-1,1\}$
 - Assume that a dummy variable 1 is included in x_n , and b is included in w
- Initialize weight vector $w = \mathbf{0}$
- Loop for T iterations (or till converge; i.e., no misclassification)
 - Loop for all training examples x_n (random order!)
 - Predict $\hat{y}_n = \text{sign}(w^T x_n)$
 - If $\hat{y}_n \neq y_n$
 - Update: $w \leftarrow w + \eta(y_n x_n)$
- Note: η is the learning rate (step size) and $\eta \geq 0$

Perceptron Algorithm

- Example:

$$w \leftarrow w + \eta(y_n x_n)$$



Perceptron Algorithm

- Why does it work?
 - If $\hat{y}_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n) = -1$ but $y_n = +1$
 - $(\mathbf{w} + \eta(y_n \mathbf{x}_n))^T \mathbf{x}_n = \mathbf{w}^T \mathbf{x}_n + \eta \textcolor{green}{y_n} \textcolor{blue}{x_n}^T \mathbf{x}_n \geq \mathbf{w}^T \mathbf{x}_n$
 - If $\hat{y}_n = \text{sign}(\mathbf{w}^T \mathbf{x}_n) = +1$ but $y_n = -1$
 - $(\mathbf{w} + \eta(y_n \mathbf{x}_n))^T \mathbf{x}_n = \mathbf{w}^T \mathbf{x}_n + \eta \textcolor{green}{y_n} \textcolor{blue}{x_n}^T \mathbf{x}_n \leq \mathbf{w}^T \mathbf{x}_n$

Perceptron Algorithm: online/streaming version

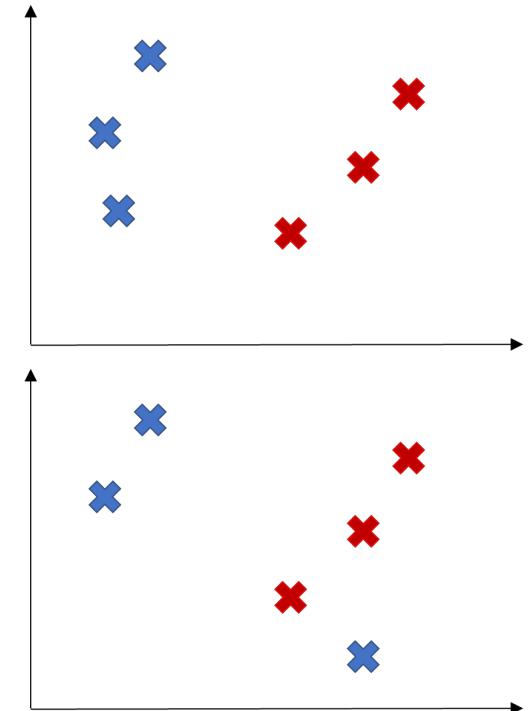
- Let $x \in \mathbb{R}^D$ and $w \in \mathbb{R}^D$, $y \in \{-1, 1\}$
 - Assume that a dummy variable 1 is included in x_n , and b is included in w
 - Initialize weight vector $w = \mathbf{0}$
 - **Online version:** Get one training data (x, y)
 - Predict $\hat{y} = \text{sign}(w^T x)$
 - If $\hat{y} \neq y$
 - **Update:** $w \leftarrow w + \eta(yx)$

Perceptron Algorithm

- Convergence (in terms of the training error):
 - If the data is linear separable

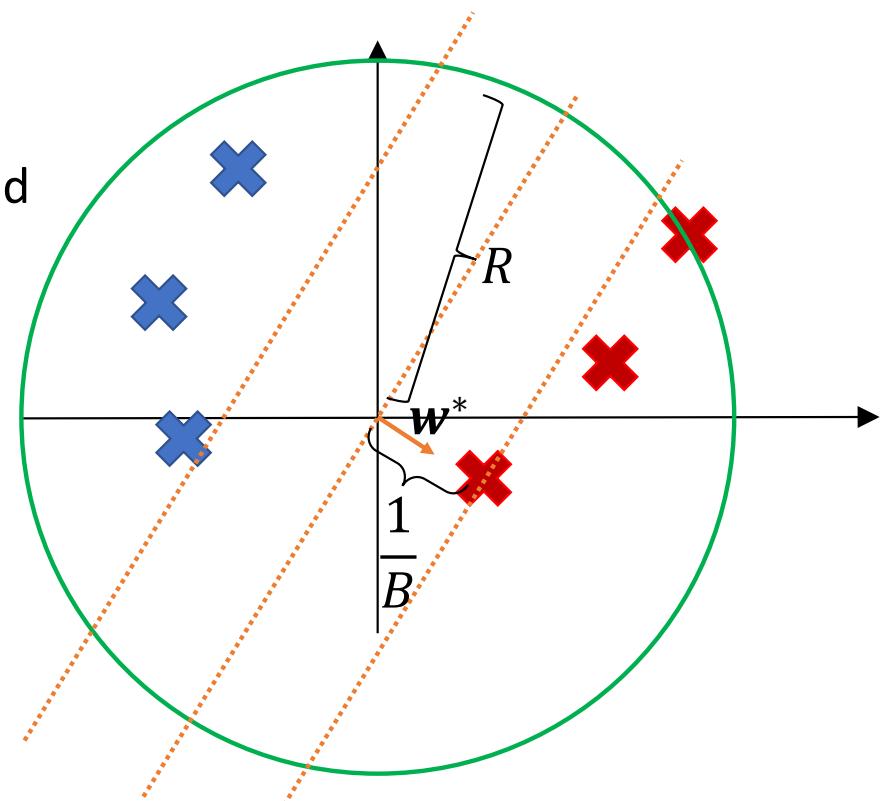
- Not convergence:
 - If the data is not linear separable
 - Not good ... any idea to resolve it?

- **The pocket algorithm:** keep the best solution so far “in the pocket”
 - “Best” in terms of the training error



Proof of convergence

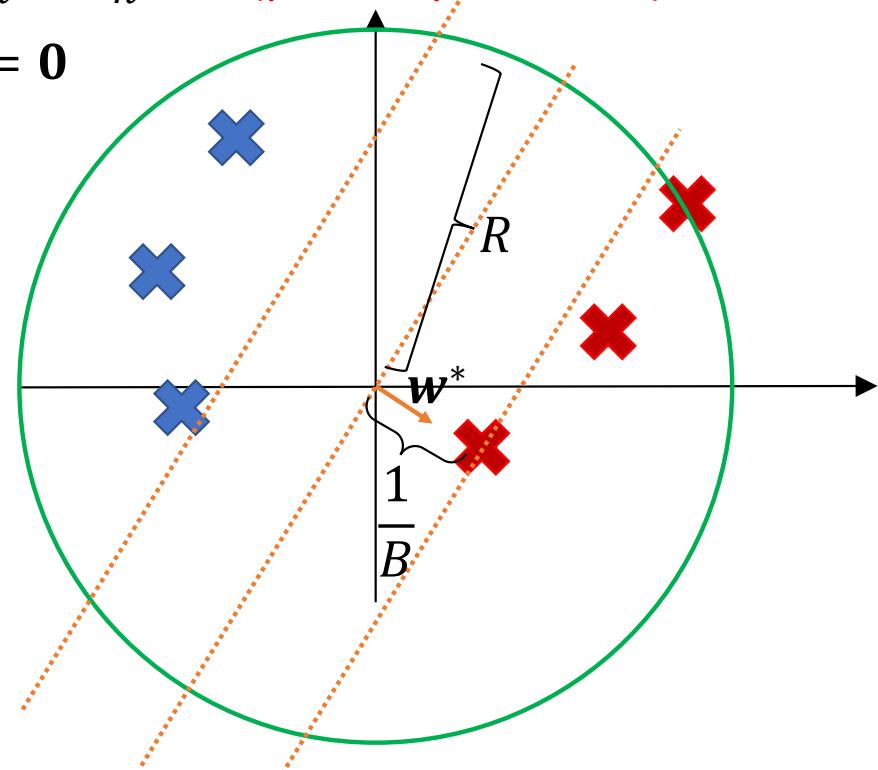
- If the training set $D_{tr} = \{(x_n \in \mathbb{R}^D, y_n \in \{-1, +1\})\}_{n=1}^N$ is linearly separable
 - $R = \max_n \|x_n\|$
 - $B = \min_n \|\mathbf{w}\|, s.t. \forall n, y_n \mathbf{w}^T \mathbf{x}_n \geq 1$
 - Stop after performing at most $(RB)^2$ updated
- Intuition:
 - The more separable the data are
 - The faster the convergence



Proof of convergence

- Proof sketch

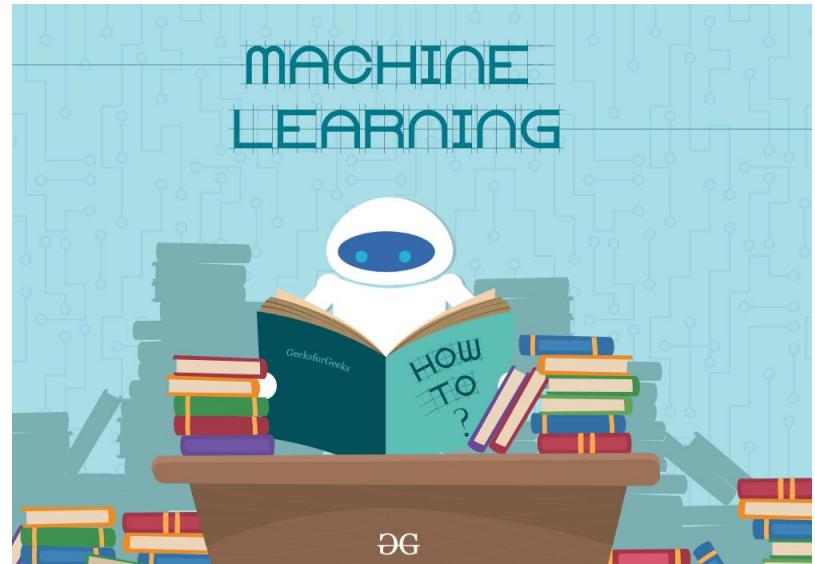
- Let \mathbf{w}^* be a solution of $\min_n \|\mathbf{w}\|_2$, s.t. $\forall n, y_n \mathbf{w}^T \mathbf{x}_n \geq 1$ (perfect prediction)
- Let $\mathbf{w}^{(t)}$ be the \mathbf{w} after t updates, and $\mathbf{w}^{(0)} = \mathbf{0}$
- Let $\eta = 1$
- Show that $\frac{\mathbf{w}^{(t)T} \mathbf{w}^*}{\|\mathbf{w}^{(t)}\|_2 \|\mathbf{w}^*\|_2} \geq \frac{\sqrt{t}}{RB}$
 - The left-hand side cannot go over 1
 - $t \leq (RB)^2$
- Intuition:
 - every update moves $\mathbf{w}^{(t)}$ closer to \mathbf{w}^*
- Details:
 - See UML: 9-9.1.2



Today – Part 2

Estimating Probability from Data

- Overview
- Case study: Bernoulli distribution
 - MLE
 - MAP
 - Bayesian estimate
- *Reading: categorical and Gaussian*
- Multi-dimensional data
 - *Reading: discrete variables*



Estimating Probability from Data

- Goal: Fit a probability model p (interchangeable with P here) to data
 - $p(X, Y)$
 - $p(Y)p(X|Y)$: generative learning
 - $p(X)p(Y|X)$: discriminative learning
 - or just $p(X)$ without labels
- Why to do so?
 - We can then use the Bayes optimal classifier $\hat{y} = \operatorname{argmax}_y P(Y = y|x)$
 - We can then perform probabilistic inference
 - We can then sample/generate new data instances
 - We can then evaluate the probability of seeing a particular data instance

Estimating Probability from Data

- Data \mathcal{D} (observation)
 - $D_{tr} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$
 - $D_{tr} = \{\mathbf{x}_n\}_{n=1}^N$
- Assumption:
 - IID (or i.i.d.): Independent and Identically Distributed
 - $p((\mathbf{x}_i, y_i), (\mathbf{x}_j, y_j)) = p_i((\mathbf{x}_i, y_i)) \times p_j((\mathbf{x}_j, y_j)) = p((\mathbf{x}_i, y_i)) \times p((\mathbf{x}_j, y_j))$
 - p_i and p_j are the same probability distribution
- Fitting to $D_{tr} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ or $D_{tr} = \{\mathbf{x}_n\}_{n=1}^N$ is conceptually similar

Estimating “Parameters of Probability” from Data

- **Maximum likelihood estimation (MLE)**

- Let $p(\mathbf{x}; \boldsymbol{\theta})$ be the probability model (for one data instance) **with parameters $\boldsymbol{\theta}$**
 - For example, a Gaussian distribution's mean and variance
- Want to find with which $\boldsymbol{\theta}$, the data D_{tr} are most likely to be observed
- $\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N; \boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} p(\mathbf{x}_1; \boldsymbol{\theta}) \times \dots p(\mathbf{x}_N; \boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \prod_n p(\mathbf{x}_n; \boldsymbol{\theta})$

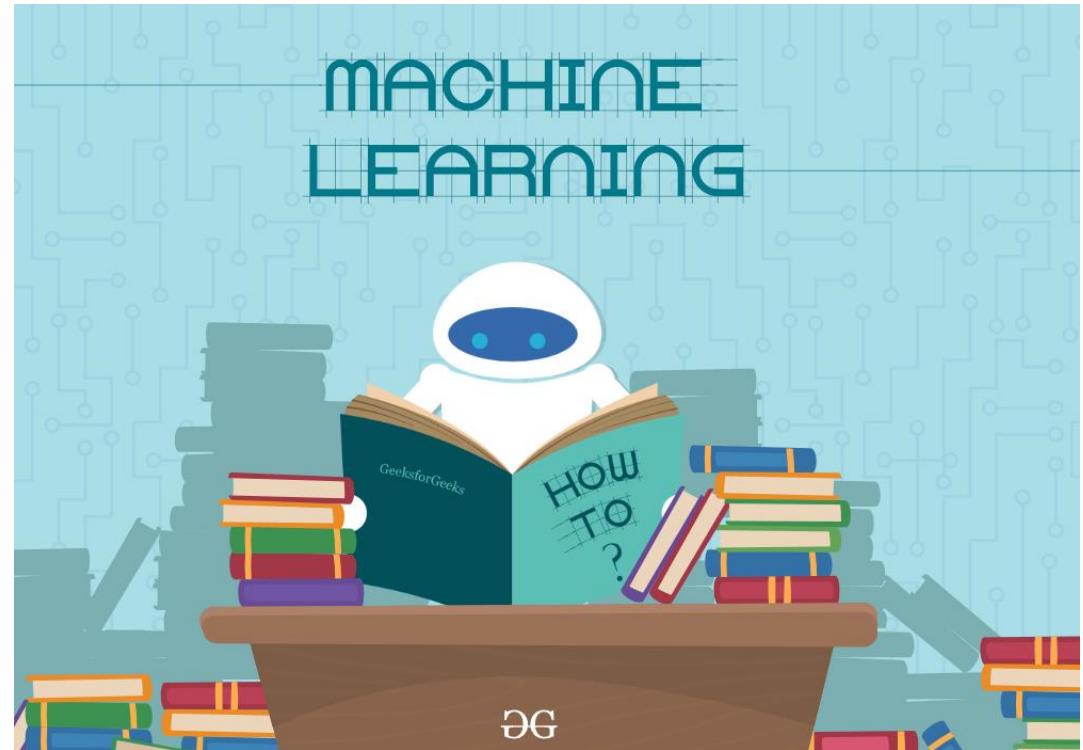
- **Maximum a posteriori estimation (MAP)**

- Assume that $\boldsymbol{\theta}$ is a random variable with a distribution $p(\boldsymbol{\theta})$; $p(\mathbf{x}; \boldsymbol{\theta})$ becomes $p(\mathbf{x}|\boldsymbol{\theta})$
 - For example, a Gaussian's mean is around 0 and variance is around 1
- Want to find $\boldsymbol{\theta}$ with a high likelihood of observing D_{tr} & a high probability $p(\boldsymbol{\theta})$
- $\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta}) \times p(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N | \boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} p(\boldsymbol{\theta}) \times \prod_n p(\mathbf{x}_n | \boldsymbol{\theta})$

Today

Estimating Probability from Data

- Overview
- Case study: Bernoulli distribution
 - MLE
 - MAP
 - Bayesian estimate
- Reading: categorical and Gaussian
- Multi-dimensional data
 - Reading: discrete variables



Bernoulli distribution (for a binary random variable)

- Parameters: $0 \leq \lambda \leq 1$

- $P(\text{"head"}) = \lambda$
 - $P(\text{"tail"}) = 1 - \lambda$

- Examples (coin flipping):

- Fair coins: $\lambda = 0.5$
 - Unfair coins: $\lambda \neq 0.5$



λ

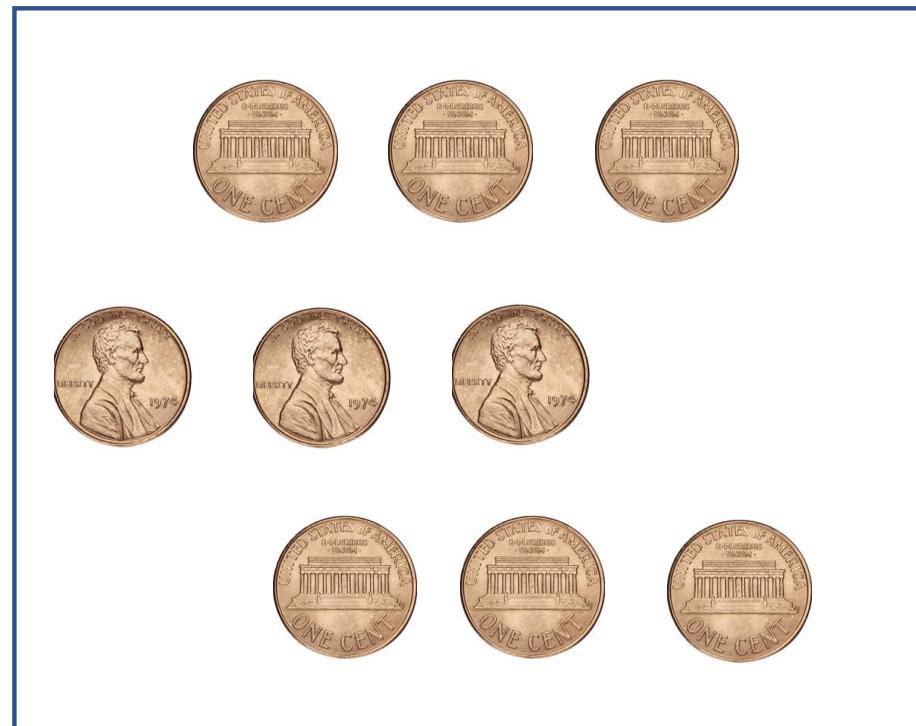
$1 - \lambda$

- Usages (given λ):

- You can tell that seeing a "head" is $P(\text{"head"}) = \lambda$
 - You can simulate the coin flipping and get a "head" with a probability λ

Bernoulli distribution: parameter estimation

- Don't know λ , but have "N" independent instances/samples
 - Ex: have 9 coin flipping samples
 - 3 heads, 6 tails
- What should λ (i.e., $P(\text{"head"})$) be?
- Why?



Bernoulli distribution: parameter estimation

- Don't know λ , but have "N" independent instances/samples
 - Ex: have 9 coin flipping samples
 - 3 heads, 6 tails
- What should λ (i.e., $P(\text{"head"})$) be?
 - Assume the first three samples are "head"
 - Assume the last six samples are "tail" ... this removes a constant $\binom{9}{3}$
 - **Joint probability:** $P(X_1 = \text{head}, X_2 = \text{head}, X_3 = \text{head}, X_4 = \text{tail}, \dots, X_9 = \text{tail}) = P(X_1 = \text{head})P(X_2 = \text{head})P(X_3 = \text{head})P(X_4 = \text{tail}) \dots P(X_9 = \text{tail}) =$ independent
 $\lambda^3(1 - \lambda)^6$
 - The "likelihood function" of λ to see these assignments/observations is $\lambda^3(1 - \lambda)^6$

Likelihood function

- In statistics, the likelihood function (simply, the likelihood) measures the goodness of fit of a model to a set of data instances for a given value of the unknown parameter λ .
- It is formed from the joint probability distribution of the samples
- It is viewed and used as a function of the parameters only, thus treating the random variables X as fixed at the observed values.
- In other words, the random variables get their assignments from the observations, and the likelihood is the joint probability of seeing these observations, given a value of the unknown parameter λ .

Assumptions on the observations

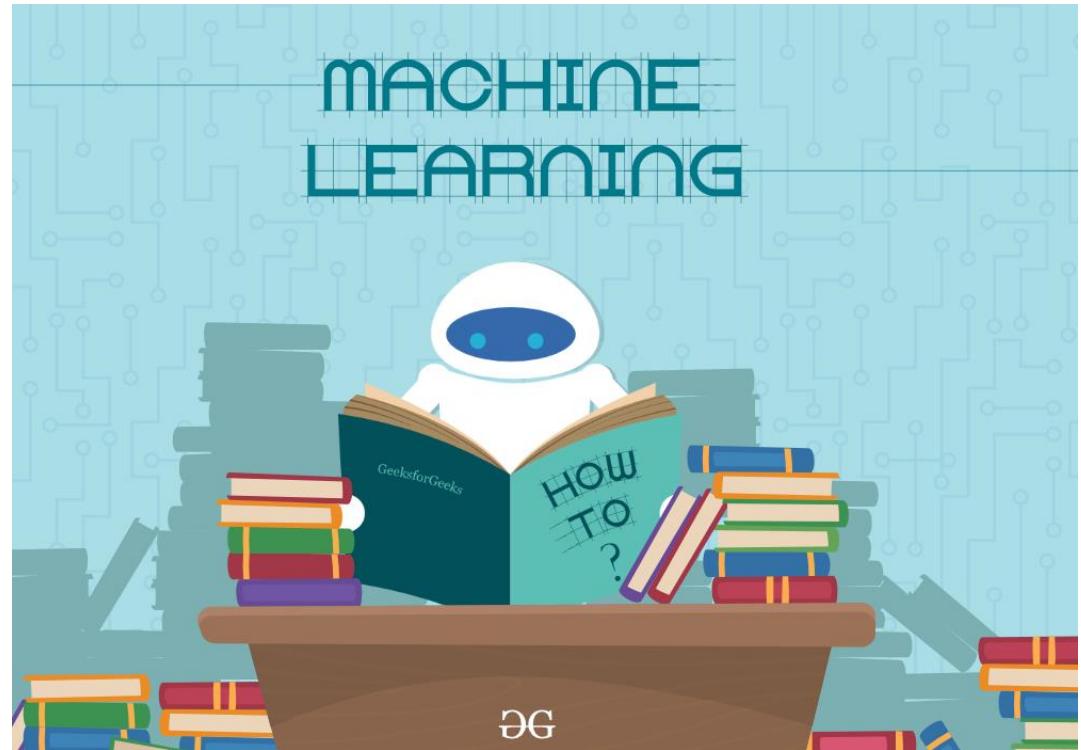
- Every coin flipping is based on the **same** probability (same parameter λ)
- Each flip is **independent**
 - Doesn't affect the outcomes of other flips
- IID (or i.i.d.): **Independent** and **Identically Distributed**



Today

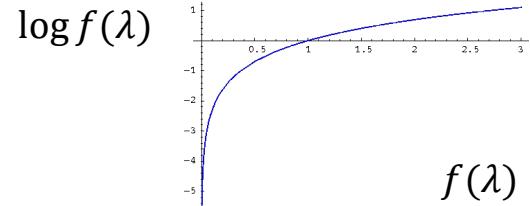
Estimating Probability from Data

- Overview
- Case study: Bernoulli distribution
 - MLE
 - MAP
 - Bayesian estimate
- Reading: categorical and Gaussian
- Multi-dimensional data
 - Reading: discrete variables



Maximum likelihood estimation (MLE)

- Likelihood function: $p(\text{observations}; \text{parameters})$
 - When IID, $p(\text{observations} ; \text{parameters}) = \prod_i p(\text{observation}_i ; \text{parameters})$
- **MLE:** choose the parameters that maximize the probability of observed data!
- Recipe:
 - Compute the log-likelihood
 - Take derivatives with respect to the parameters
 - Equate to zero and solve
 - You should check if you find a minimum or a maximum, by checking the second-order derivative (i.e., Hessian)



Maximum likelihood estimation (MLE)

- What $0 \leq \lambda \leq 1$ will maximize $\lambda^3(1 - \lambda)^6$?

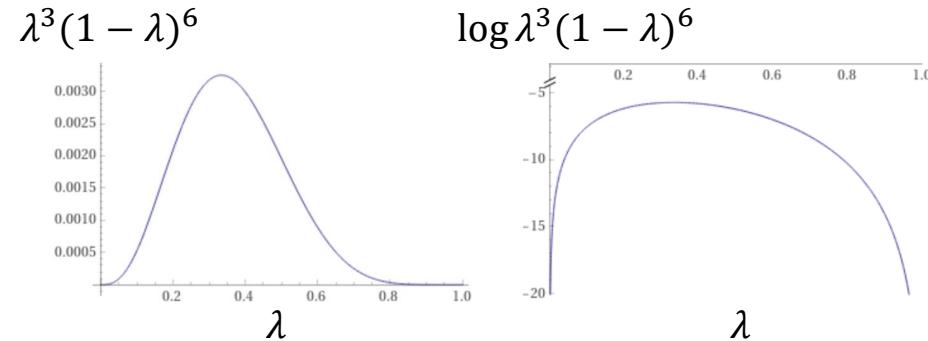
- Log: monotonically increasing
- Turn a multiplication into sum

$$\begin{aligned}\log \lambda^3(1 - \lambda)^6 &= \log \lambda^3 + \log(1 - \lambda)^6 \\ &= 3 \log \lambda + 6 \log(1 - \lambda)\end{aligned}$$

- Derivative = 0:

$$\frac{d \log \lambda^3(1-\lambda)^6}{d\lambda} = 3 \frac{d \log \lambda}{d\lambda} + 6 \frac{d \log(1-\lambda)}{d\lambda} = 3 \frac{1}{\lambda} - 6 \frac{1}{1-\lambda} = \frac{3-9\lambda}{\lambda(1-\lambda)} = 0$$

- Answer: $\lambda = \frac{3}{9} = \frac{\text{\# of heads}}{\text{\# of samples}}$

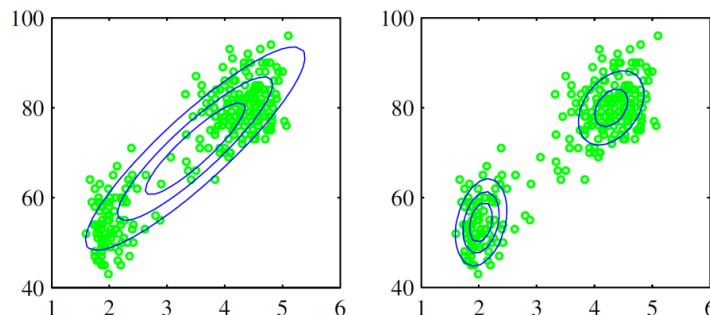


Maximum likelihood estimation (MLE)

- Let X_i ("head") = 1, X_i ("tail") = 0
 - $P(X_1 = x_1, \dots, X_N = x_N) = P(X_1 = x_1) \dots P(X_N = x_N) = \prod_{i=1}^N P(X_i = x_i)$
 - $\prod_{i=1}^N P(X_i = x_i) = \prod_{i=1}^N \lambda^{x_i} (1 - \lambda)^{1-x_i}$
 - Ex. $x_i = 1, \lambda^{x_i} (1 - \lambda)^{1-x_i} = \lambda^1 (1 - \lambda)^{1-1} = \lambda$
 - EX. $x_i = 0, \lambda^{x_i} (1 - \lambda)^{1-x_i} = \lambda^0 (1 - \lambda)^{1-0} = 1 - \lambda$
 - $\log \prod_{i=1}^N \lambda^{x_i} (1 - \lambda)^{1-x_i} = \sum_{i=1}^N \log \{\lambda^{x_i} (1 - \lambda)^{1-x_i}\} = \sum_{i=1}^N \{\log \lambda^{x_i} + \log (1 - \lambda)^{1-x_i}\}$
 $= \sum_{i=1}^N x_i \log \lambda + \sum_{i=1}^N (1 - x_i) \log (1 - \lambda)$
 $= (\#x_i = 1) \log \lambda + (\#x_i \neq 1) \log (1 - \lambda)$
 - Maximum likelihood estimation: $\underset{\lambda}{\operatorname{argmax}} P(X_1 = x_1, \dots, X_N = x_N; \lambda) = \underset{\lambda}{\operatorname{argmax}} P(\mathcal{D}; \lambda)$
 - $\lambda = \frac{\text{\# of 1}}{\text{\# of samples}}$

Quick summary of MLE

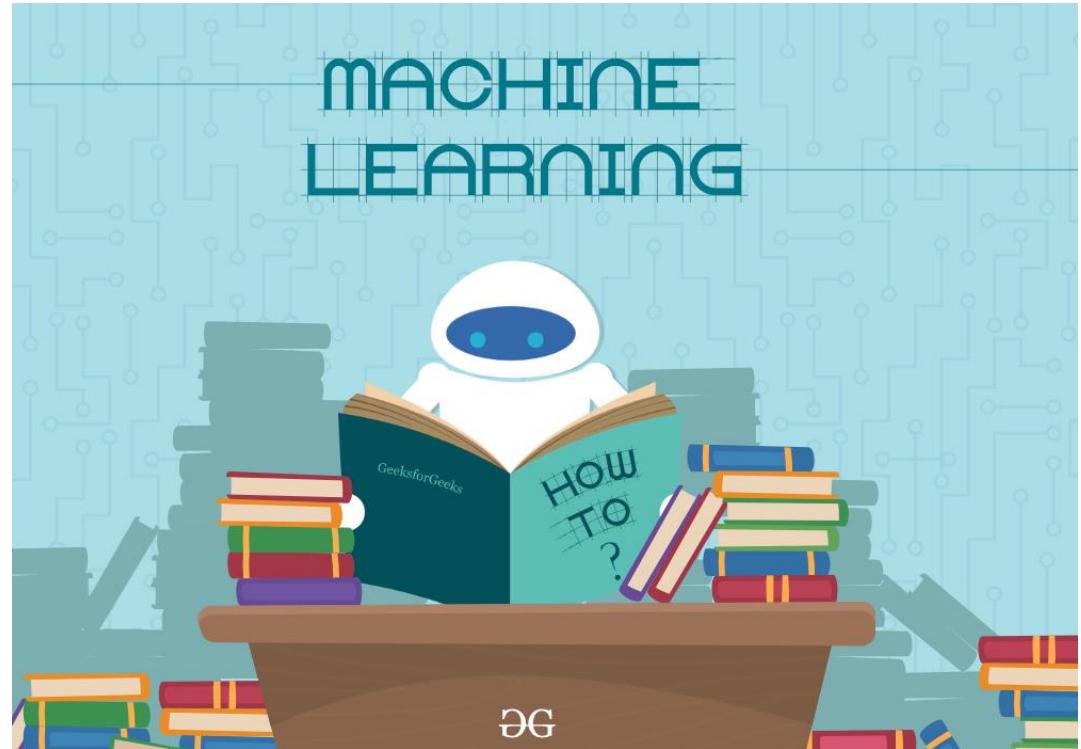
- MLE gives the explanation of the data you observed
- MLE finds the **TRUE** parameters if
 - You choose the right probability model (hypothesis class \mathcal{H})
 - N is large enough (seen from the concentration theory, Hoeffding's inequality)
- When N is small, MLE can overfit the data
- When the model is wrong, whatever you estimate is likely meaningless ...



Today

Estimating Probability from Data

- Overview
- Case study: Bernoulli distribution
 - MLE
 - MAP
 - Bayesian estimate
- Reading: categorical and Gaussian
- Multi-dimensional data
 - Reading: discrete variables



How to prevent overfitting

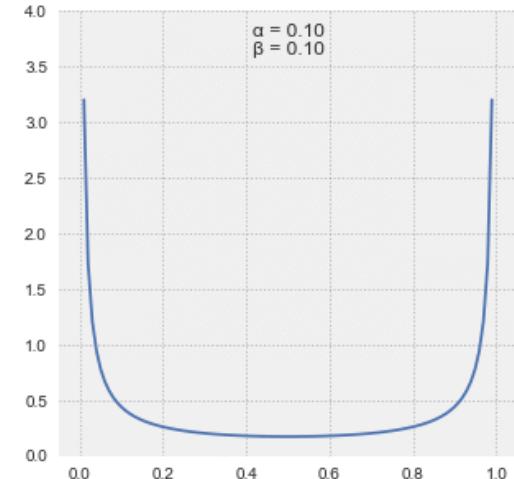
- Overfitting:
 - You model fits your training data too well (even the noise in the data)
 - N is small
 - The hypothesis class is too complicated
 - You model does not work well for your test data
- Here, to resolve the issue when N is small
 - Adding some **prior belief** “before” you see the data
 - Ex: λ for $P(\text{“head”}) = \lambda$ is close 0.5

Maximum a posteriori estimation (MAP)

- Let λ be a random variable as well, with a probability $P(\lambda)$
 - Ex. You have some prior knowledge about the coin!
- $\operatorname{argmax}_{\lambda} P(\lambda|\mathcal{D}) = \operatorname{argmax}_{\lambda} \frac{P(\mathcal{D}|\lambda)p(\lambda)}{P(\mathcal{D})} = \operatorname{argmax}_{\lambda} P(\mathcal{D}|\lambda)p(\lambda)$ Bayes' rule!
 - $p(\lambda)$: prior distribution over λ , before seeing \mathcal{D}
 - $P(\mathcal{D}|\lambda)$: likelihood of the data given λ
 - $P(\lambda|\mathcal{D})$: posterior distribution over λ , after seeing \mathcal{D}
- How to set $P(\lambda)$, the prior probability of λ , depends on:
 - Domain knowledge; or learned (e.g., empirical Bayes)
 - “Specific” probability models (guess what it will be?)

Maximum a posteriori estimation (MAP)

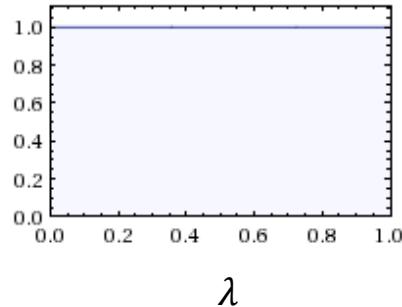
- Let $P(\lambda) = \frac{\lambda^{\alpha-1} (1-\lambda)^{\beta-1}}{Z(\alpha, \beta)}$
 - This is called the Beta distribution of λ
 - α and β are the parameters of the distribution
 - $Z(\alpha, \beta)$ is a constant (given α and β) for normalization
- Why do we use the Beta distribution?
 - It models the probability of λ , where λ has a range = [0, 1]
 - Mathematically beautiful
 - $P(\lambda|\mathcal{D}) \propto P(\mathcal{D}|\lambda)p(\lambda) \propto (\lambda^{\# \text{ of head}} (1 - \lambda)^{\# \text{ of tail}})(\lambda^{\alpha-1} (1 - \lambda)^{\beta-1}) \propto \lambda^{\# \text{ of head} + \alpha - 1} (1 - \lambda)^{\# \text{ of tail} + \beta - 1}$
 - Together with Bernoulli (or precisely binomial) as the likelihood, the posterior is again Beta
 - Such a prior is called the **conjugate prior** (of a certain likelihood function)!



Maximum a posteriori estimation (MAP)

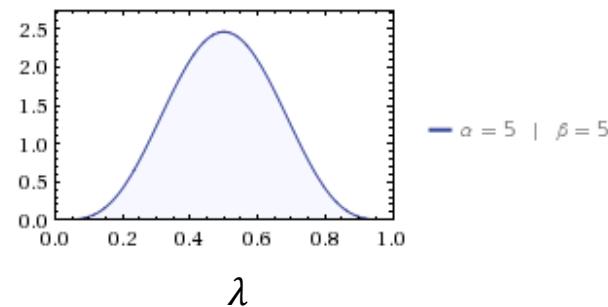
Beta(1,1) (λ)

Plots:



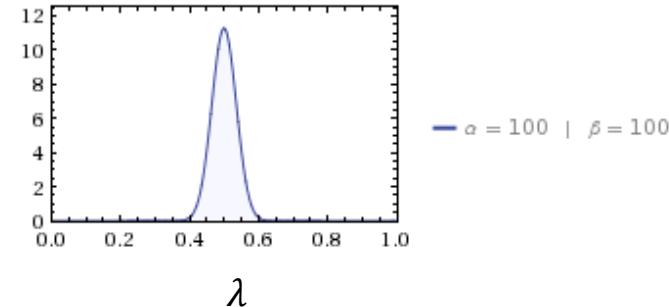
Beta(5,5) (λ)

Plots:

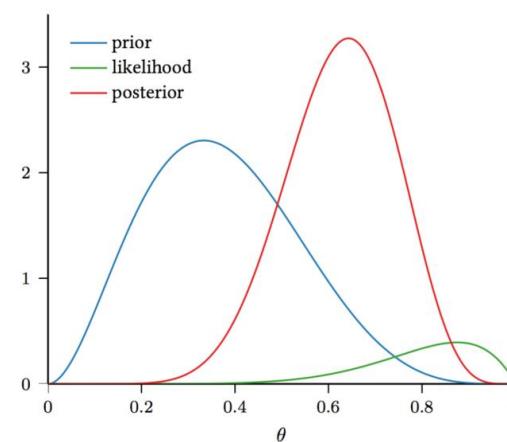
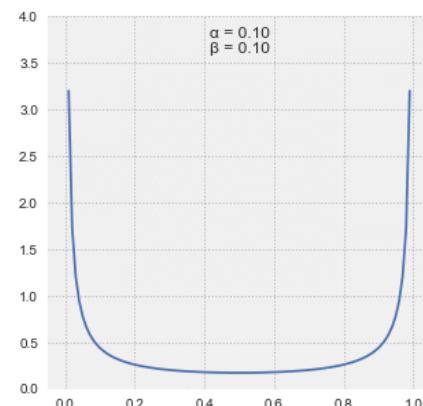


Beta(100,100) (λ)

Plots:



$$P(\lambda) = \frac{\lambda^{\alpha-1} (1-\lambda)^{\beta-1}}{Z(\alpha, \beta)}$$



$$P(\lambda|D) \propto P(D|\lambda)p(\lambda)$$

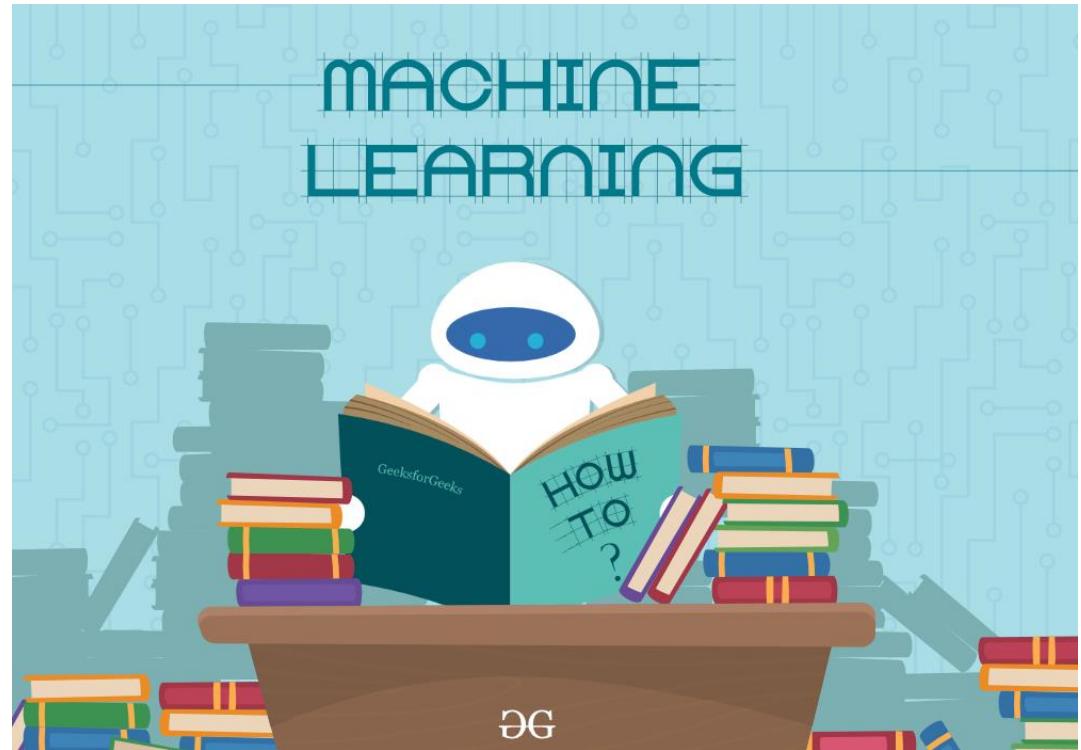
Maximum a posteriori estimation (MAP)

- $p(\mathcal{D}|\lambda)p(\lambda) \propto \lambda^{\# \text{ of heads}} \times (1 - \lambda)^{\# \text{ of tails}} \times \lambda^{\alpha-1} \times (1 - \lambda)^{\beta-1}$
 $= \lambda^{\# \text{ of heads} + \alpha - 1} \times (1 - \lambda)^{\# \text{ of tails} + \beta - 1}$
- What is $\underset{\lambda}{\operatorname{argmax}} P(\lambda|\mathcal{D}) = \underset{\lambda}{\operatorname{argmax}} \log p(\mathcal{D}|\lambda) + \log p(\lambda)$?
 - $\lambda = \frac{\# \text{ of heads} + \alpha - 1}{\# \text{ of samples} + \alpha + \beta - 2}$
- $\alpha - 1, \beta - 1$ (or α, β) are called pseudo counts
 - When N is large, MAP is similar to MLE

Today

Estimating Probability from Data

- Overview
- Case study: Bernoulli distribution
 - MLE
 - MAP
 - Bayesian estimate
- Reading: categorical and Gaussian
- Multi-dimensional data
 - Reading: discrete variables



Comparison: MLE, MAP, and Bayesian

- Given “NH heads and NT tails”
 - MLE: $\underset{\lambda}{\operatorname{argmax}} P(\mathcal{D}|\lambda) = \frac{NH}{NH+NT}$
 - MAP: $\underset{\lambda}{\operatorname{argmax}} P(\lambda|\mathcal{D}) = \underset{\lambda}{\operatorname{argmax}} P(\mathcal{D}|\lambda)P(\lambda) = \frac{NH+\alpha-1}{(NH+\alpha-1)+(NT+\beta-1)}$, if $P(\lambda)$ is a Beta Dist.
 - Bayesian estimates: just **keep $P(\lambda|\mathcal{D})$** , with λ still a random variable
 - This keeps more information about λ (not just the mode)
 - MLE and MAP are point estimates: choosing one λ ; Bayesian: keep a distribution
- Usages after seeing \mathcal{D} : future prediction X = “head”
 - $P(X = \text{"head"}; \lambda) = \lambda$ λ is estimated by MLE or MAP
 - $P(X = \text{"head"}|\mathcal{D}) = \int P(X = \text{"head"}|\lambda) P(\lambda|\mathcal{D}) d\lambda = \frac{NH+\alpha}{(NH+\alpha)+(NT+\beta)}$; “Bayesian inference”

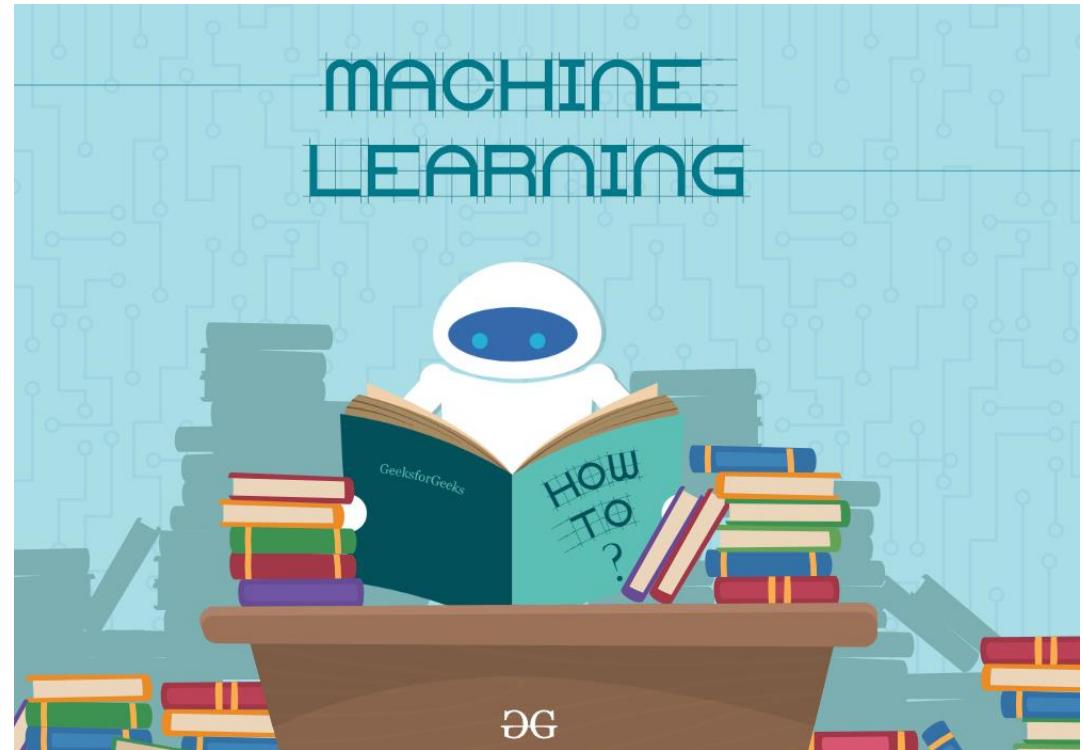
More about Bayesian

- $P(X = \text{"head"} | \mathcal{D}) = \int P(X = \text{"head"}, \lambda | D) d\lambda =$
- $\int P(X = \text{"head"} | \lambda, D) P(\lambda | \mathcal{D}) d\lambda = \int P(X = \text{"head"} | \lambda) P(\lambda | \mathcal{D}) d\lambda =$
- $\int \lambda P(\lambda | \mathcal{D}) d\lambda = E[\lambda | \mathcal{D}] = \frac{NH + \alpha}{(NH + \alpha) + (NT + \beta)}$
- In general (for other distributions), the integration is intractable, and we must approximate it by Monte Carlo methods, variational methods, etc.
 - $\int P(X=x | \lambda) P(\lambda | \mathcal{D}) d\lambda \approx \frac{1}{M} \sum_m P(X=x | \lambda_m)$, where $\lambda_m \sim P(\lambda | \mathcal{D})$
- $P(\lambda | \mathcal{D}) = \frac{P(\lambda, \mathcal{D})}{P(\mathcal{D})}$ is usually intractable, due to $P(\mathcal{D}) = \int P(\mathcal{D} | \lambda) P(\lambda) d\lambda$

Today

Estimating Probability from Data

- Overview
- Case study: Bernoulli distribution
 - MLE
 - MAP
 - Bayesian estimate
- **Reading: categorical and Gaussian**
- Multi-dimensional data
 - Reading: discrete variables



Reading: Categorical Distribution



THE OHIO STATE UNIVERSITY

Categorical distribution

- Parameters: $0 \leq \lambda[c] \leq 1, \sum_{c=1}^C \lambda[c] = 1$

- $P("1") = \lambda[1]$

- $P("6") = \lambda[6] = 1 - \sum_{c=1}^5 \lambda[c]$



- Examples (dice tossing):

- Fair dice: $\lambda[c] = 1/6$

- Unfair dice: $\exists c, \lambda[c] \neq 1/6$

- Usages (given $\lambda[c]$):

- The probability to see a “1” is $P("1") = \lambda[1]$

- You can toss a dice to get a “1” with a probability $\lambda[1]$

Categorical distribution: parameter estimation

- Don't know $\lambda[c]$, but have "N" independent instances/samples
 - Ex: have 100 samples of tossing a dice
 - 30 of them are "1", 25 are "2", 20 are "3", 15 are "4", 5 are "5", and 5 are "6"
- What should $\lambda[c]$ be for each c ?
- Why?

Categorical distribution: parameter estimation

- Don't know $\lambda[c]$, but have "N" independent instances/samples
 - Ex: have 100 samples of tossing a dice
 - 30 of them are "1", 25 are "2", 20 are "3", 15 are "4", 5 are "5", and 5 are "6"
- What should $\lambda[c]$ be for each c ?
 - Assume the first 30 samples are "1", and so on
 - $P(X_1 = "1", \dots, X_{100} = "6") = P(X_1 = "1") \dots P(X_{100} = "6") = \lambda[1]^{30} \lambda[2]^{25} \lambda[3]^{20} \lambda[4]^{15} \lambda[5]^5 \lambda[6]^5$
 - The "likelihood" to see these assignments is $\lambda[1]^{30} \lambda[2]^{25} \lambda[3]^{20} \lambda[4]^{15} \lambda[5]^5 \lambda[6]^5$

Maximum likelihood estimation (MLE)

- What $0 \leq \lambda[c] \leq 1$ will maximize the likelihood?
 - Log: monotonically increasing

$$\begin{aligned} & \log \lambda[1]^{30} \lambda[2]^{25} \lambda[3]^{20} \lambda[4]^{15} \lambda[5]^5 \lambda[6]^5 \\ &= 30 \log \lambda[1] + 25 \log \lambda[2] + \dots + 5 \log \lambda[6] \end{aligned}$$

- Derivative = 0 (with a constraint $\sum_{c=1}^C \lambda[c] = 1$, by Lagrangian multipliers):

$$\frac{d(30 \log \lambda[1] + \dots + 5 \log \lambda[6] + \gamma(\sum_{c=1}^C \lambda[c] - 1))}{d\lambda[c]} = (\# \text{ of } c) \frac{1}{\lambda[c]} + \gamma = 0$$

- Answer: $\lambda[c] = \frac{\# \text{ of } c}{\# \text{ of samples}}$

Maximum likelihood estimation (MLE)

- Let X ("1") = 1, X ("6") = 6
 - $P(X_1 = x_1, \dots, X_N = x_N) = P(X_1 = x_1) \dots P(X_N = x_N) = \prod_{i=1}^N P(X_i = x_i)$
 - $\prod_{i=1}^N P(X_i = x_i) = \prod_{i=1}^N \left\{ \prod_{c=1}^C \lambda[c] \mathbf{1}[x_i=c] \right\}$
 - Ex. $x_i = 1, \prod_{c=1}^C \lambda[c] \mathbf{1}[x_i=c] = \lambda[1]$
 - EX. $x_i = 6, \prod_{c=1}^C \lambda[c] \mathbf{1}[x_i=c] = \lambda[6]$

Indicator function

Maximum likelihood estimation (MLE)

- Let X ("1") = 1, X ("6") = 6
 - $P(X_1 = x_1, \dots, X_N = x_N) = P(X_1 = x_1) \dots P(X_N = x_N) = \prod_{i=1}^N P(X_i = x_i)$
 - $\prod_{i=1}^N P(X_i = x_i) = \prod_{i=1}^N \left\{ \prod_{c=1}^C \lambda[c]^{1[x_i=c]} \right\}$
 - Ex. $x_i = 1, \prod_{c=1}^C \lambda[c]^{1[x_i=c]} = \lambda[1]$
 - EX. $x_i = 6, \prod_{c=1}^C \lambda[c]^{1[x_i=c]} = \lambda[6]$
 - $\log \prod_{i=1}^N \prod_{c=1}^C \lambda[c]^{1[x_i=c]} = \sum_{i=1}^N \sum_{c=1}^C \mathbf{1}[x_i=c] \log \lambda[c]$
 $= \sum_{c=1}^C \log \lambda[c] \sum_{i=1}^N \mathbf{1}[x_i=c] = \sum_{c=1}^C \log \lambda[c] (\#x_i = c)$
 - Maximum likelihood estimation: $\underset{\lambda}{\operatorname{argmax}} P(X_1 = x_1, \dots, X_N = x_N | \lambda) = \underset{\lambda}{\operatorname{argmax}} P(\mathcal{D} | \lambda)$
 - $\lambda[c] = \frac{\text{\# of } c}{\text{\# of samples}}$

Maximum a posteriori estimation (MAP)

- Let λ be a random vector as well, with a probability $P(\lambda)$
 - Ex. You have some prior knowledge about the dice!
- $\operatorname{argmax}_{\lambda} P(\lambda|\mathcal{D}) = \operatorname{argmax}_{\lambda} \frac{P(\mathcal{D}|\lambda)p(\lambda)}{P(\mathcal{D})} = \operatorname{argmax}_{\lambda} P(\mathcal{D}|\lambda)p(\lambda)$ Bayes' rule!
- How to set $P(\lambda)$, the prior probability of λ , depends on:
 - Domain knowledge
 - “Specific” probability models (guess what it will be?)

Summary: categorical distribution

- Parameters: $0 \leq \lambda[c] \leq 1, \sum_{c=1}^C \lambda[c] = 1$

- $\circ P("1") = \lambda[1]$

- $\circ P("6") = \lambda[6] = 1 - \sum_{c=1}^5 \lambda[c]$



- MLE: $\hat{\lambda}[c] = \frac{\text{\# of } c}{\text{\# of samples}}$

- \circ Including a Lagrangian multipliers $\gamma(\sum_{c=1}^C \lambda[c] - 1)$ in optimization

- MAP:

- \circ Dirichlet distribution: $p(\boldsymbol{\lambda}) = \frac{\prod_c \lambda[c]^{\alpha[c]-1}}{Z(\boldsymbol{\alpha})}$

- $\circ \lambda[c] = \frac{\text{\# of } c + (\alpha[c]-1)}{\text{\# of samples} + \sum_{c=1}^C \alpha[c]-1}$

Reading: Gaussian



THE OHIO STATE UNIVERSITY

Gaussian distribution: continuous variables

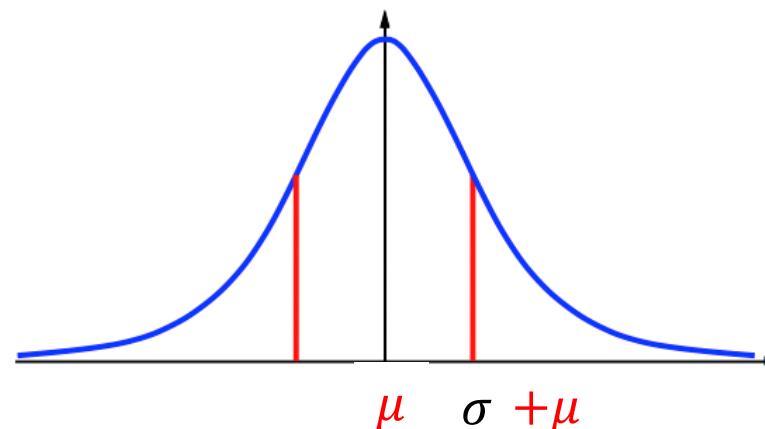
- Let X be a random variable for coin weight, what is $P(X = x)$?
- Gaussian is a popular and widely applicable distribution for continuous R.V.

$$P(X = x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$$

- Properties: uni-modal

- Parameters:

- Mean μ
- Variance σ^2



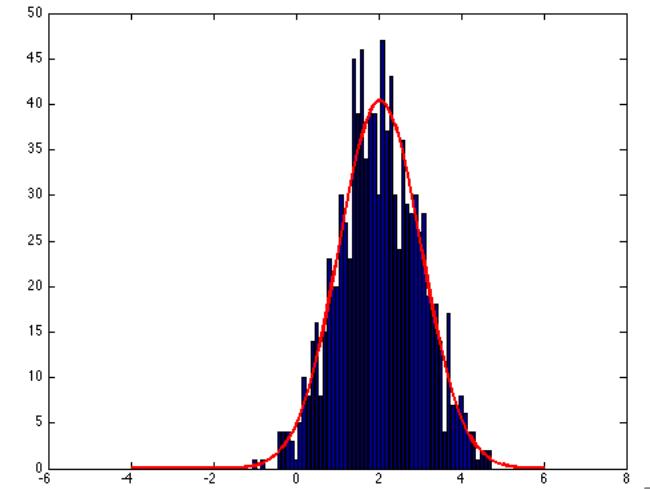
Gaussian distribution: continuous variables

- Don't know μ and σ , but have "N" independent instances/samples
 - Ex: have 6 coins
 - Weight: 4.0, 5.2, 5.3, 5.5, 6.0, 4.8
 - Gaussian is suitable for some "uni-modal" data distributions, but not for multi-modal
- What should μ and σ be?
 - $P(X_1 = 4.0, X_2 = 5.2, \dots, X_6 = 4.8) = P(X_1 = 4.0)P(X_2 = 5.2) \dots P(X_6 = 4.8) = \prod_{i=1}^N P(X_i = x_i) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x_i - \mu)^2 / 2\sigma^2}$

Maximum likelihood estimation (MLE)

- What μ and σ will maximize the likelihood?
 - Log: monotonically increasing

$$\begin{aligned} \log \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x_i - \mu)^2/2\sigma^2} &= \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x_i - \mu)^2/2\sigma^2} \right) \\ &= \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) + \log \left(e^{-(x_i - \mu)^2/2\sigma^2} \right) \\ &= \sum_{i=1}^N -\log \left(\sqrt{2\pi\sigma^2} \right) - (x_i - \mu)^2/2\sigma^2 \end{aligned}$$



Maximum likelihood estimation (MLE)

- What μ and σ will maximize the likelihood?
 - Derivative = 0 (for μ):

$$\begin{aligned} & \frac{d \left[\sum_{i=1}^N -\log(\sqrt{2\pi\sigma^2}) - (x_i - \mu)^2 / 2\sigma^2 \right]}{d\mu} \\ &= \frac{d \left[\sum_{i=1}^N -(x_i - \mu)^2 / 2\sigma^2 \right]}{d\mu} \\ &= \frac{\sum_{i=1}^N -2(x_i - \mu)(-1)}{2\sigma^2} = \frac{\sum_{i=1}^N (x_i - \mu)}{\sigma^2} = \frac{(\sum_{i=1}^N x_i) - (N\mu)}{\sigma^2} = 0 \\ & \quad \sum_{i=1}^N x_i = N\mu \end{aligned}$$

- Answer: $\mu = \frac{\sum_{i=1}^N x_i}{N} = \text{mean}$

Maximum likelihood estimation (MLE)

- What μ and σ will maximize the likelihood?
 - Derivative = 0 (for σ):

$$\frac{d \left[\sum_{i=1}^N -\log(\sqrt{2\pi\sigma^2}) - (x_i - \mu)^2 / 2\sigma^2 \right]}{d\sigma}$$

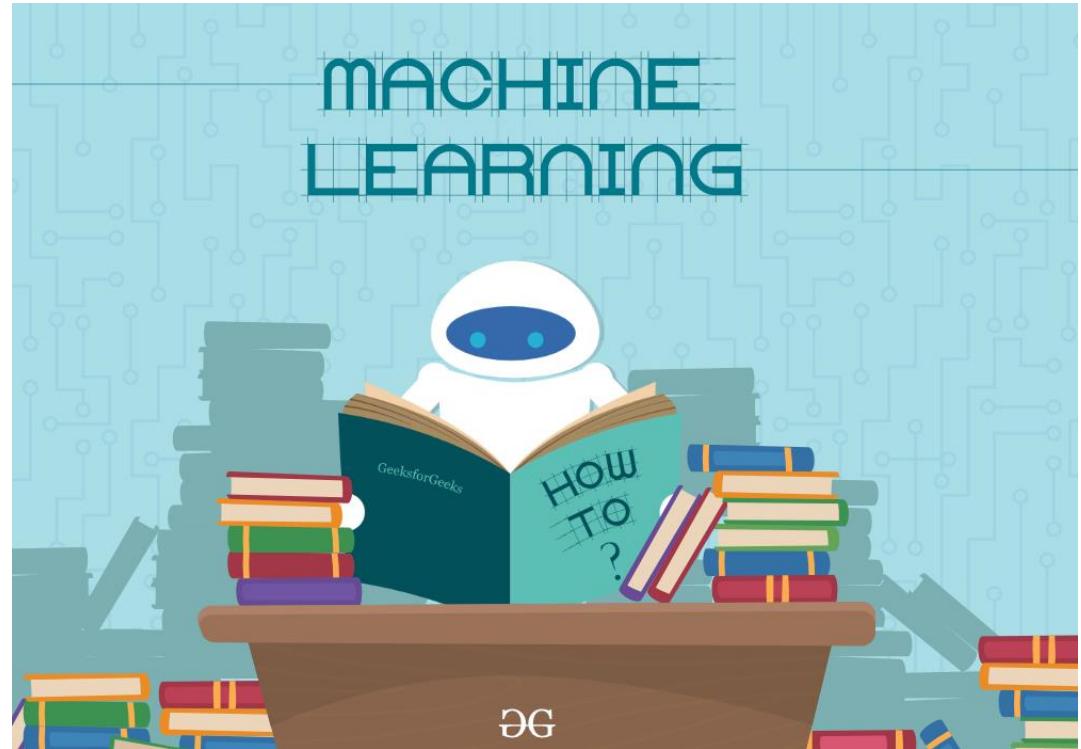
$$\begin{aligned} &= -\frac{d \left[\sum_{i=1}^N \log(\sqrt{2\pi\sigma^2}) \right]}{d\sigma} - \frac{d \left[\sum_{i=1}^N \frac{(x_i - \mu)^2}{2\sigma} \right]}{d\sigma} = \frac{dN\log(\sqrt{2\pi}\sigma)}{d\sigma} - \frac{d \left[\sum_{i=1}^N (x_i - \mu)^2 / 2\sigma^2 \right]}{d\sigma} \\ &= -\frac{N \times \sqrt{2\pi}}{\sqrt{2\pi}\sigma} - \sum_{i=1}^N -\frac{(x_i - \mu)^2}{\sigma^3} = -\frac{N\sigma^2}{\sigma^3} + \sum_{i=1}^N \frac{(x_i - \mu)^2}{\sigma^3} = 0 \end{aligned}$$

○ Answer: $\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N} = \text{variance}$

Today

Estimating Probability from Data

- Overview
- Case study: Bernoulli distribution
 - MLE
 - MAP
 - Bayesian estimate
- Reading: categorical and Gaussian
- Multi-dimensional data
 - Reading: discrete variables



Reading: Multi-dimensional discrete



THE OHIO STATE UNIVERSITY

For multi-dimensional data: discrete

- Example:

- Every time, sample one coin
- Flip the picked coin



| Instance ID: i | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|-------------------|---|---|---|---|---|---|---|---|
| Coin | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| Head (1)/Tail (0) | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 0 |

x_i
 $x_i[1]$
 $x_i[2]$

- What parameters?

- How many “joint outcomes”?

For multi-dimensional data: discrete

- Example:

- Every time, sample one coin
- Flip the picked coin



| Instance ID: i | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|-------------------|---|---|---|---|---|---|---|---|
| Coin | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| Head (1)/Tail (0) | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 0 |

x_i
 $x_i[1]$
 $x_i[2]$

- What parameters?

- How many “joint outcomes”?
- Need a parameter for each joint outcome (categorical): $\lambda[0, 0], \lambda[0, 1], \lambda[1, 0], \lambda[1, 1]$
 - $P(\text{"coin} = 0", \text{"tail"})$, $P(\text{"coin} = 0", \text{"head"})$, $P(\text{"coin} = 1", \text{"tail"})$, $P(\text{"coin} = 1", \text{"head"})$
 - MLE: counts and normalize: $P(\text{"coin} = 0", \text{"tail"}) = \frac{\# \text{ of } (\text{"coin} = 0", \text{"tail"})}{\# \text{ of samples}}$

Exponentially growing!

For multi-dimensional data: discrete

- Example:

- Every time, sample one coin
- Flip the picked coin



| Instance ID: i | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|-------------------|---|---|---|---|---|---|---|---|
| Coin | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| Head (1)/Tail (0) | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 0 |

x_i
 $x_i[1]$
 $x_i[2]$

- MLE practice

- $P(\text{"coin} = 0")?$
- $P(\text{"tail"} \mid \text{"coin} = 0")?$

$$\begin{aligned} \blacksquare \quad \frac{P(\text{"tail"}, \text{"coin} = 0)}{P(\text{"coin} = 0)} &= \frac{P(\text{"tail"}, \text{"coin} = 0)}{P(\text{"tail"}, \text{"coin} = 0) + P(\text{"head"}, \text{"coin} = 0)} = \frac{\# \text{ of } (\text{"coin} = 0, \text{"tail"})}{\# \text{ "coin} = 0} \\ \blacksquare \quad \text{Look at columns of "coin} = 0", \text{ among which how many of them are tails?} \end{aligned}$$

CSE 5523: Generative vs. Discriminative Learning



THE OHIO STATE UNIVERSITY

HW

- HW-1 is due today at midnight (11:59 PM ET).
- HW-2 will be released today by midnight, and due in two weeks (10/1).

Today

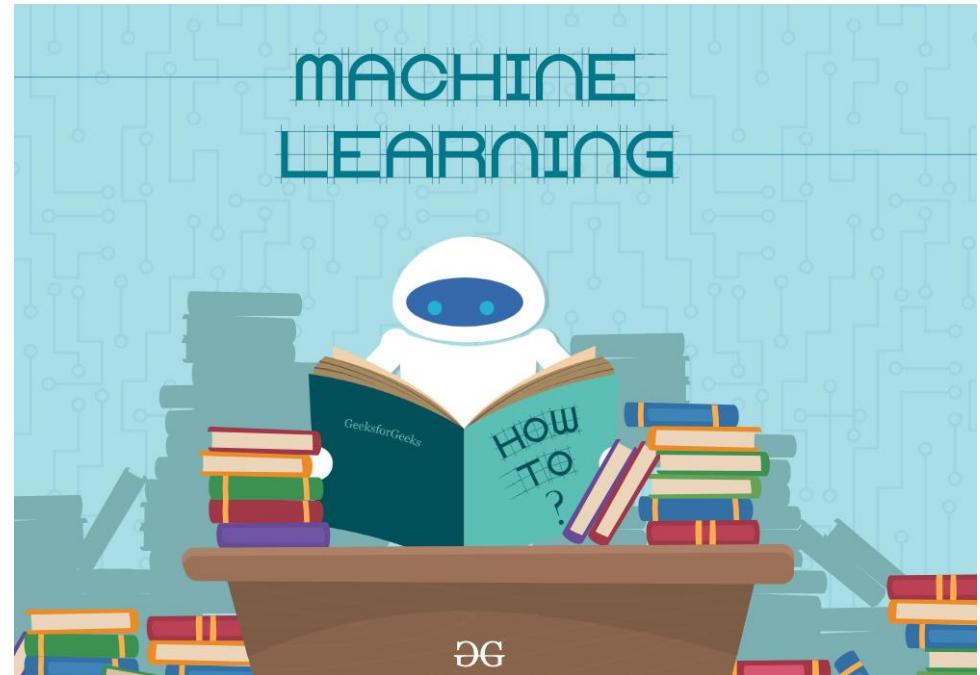
Review

Multi-dimensional Gaussian

Generative vs. discriminative learning

Linear regression revisit

Regularization, cross-validation



Estimating Probability from Data

- Goal: Fit a probability model p (interchangeable with P here) to data
 - $p(X, Y)$
 - $p(Y)p(X|Y)$: generative learning
 - $p(X)p(Y|X)$: discriminative learning
 - or just $p(X)$ without labels
- Data \mathcal{D} (observation/assignments): $D_{tr} = \{(x_n, y_n)\}_{n=1}^N$ or $D_{tr} = \{x_n\}_{n=1}^N$
 - Treatments are conceptually the same
- Assumption: **IID** (or i.i.d.): **Independent** and **Identically Distributed**

Estimating “Parameters of Probability” from Data

- Maximum likelihood estimation (MLE)

- $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} p(\mathcal{D}; \theta) = \underset{\theta}{\operatorname{argmax}} p(x_1, x_2, \dots, x_N; \theta) = \underset{\theta}{\operatorname{argmax}} p(x_1; \theta) \times \dots \times p(x_N; \theta) = \underset{\theta}{\operatorname{argmax}} \prod_n p(x_n; \theta)$

- Maximum a posteriori estimation (MAP)

- $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} P(\theta | \mathcal{D}) = \underset{\theta}{\operatorname{argmax}} p(\theta) \times p(x_1, x_2, \dots, x_N | \theta) = \underset{\theta}{\operatorname{argmax}} p(\theta) \times \prod_n p(x_n | \theta)$

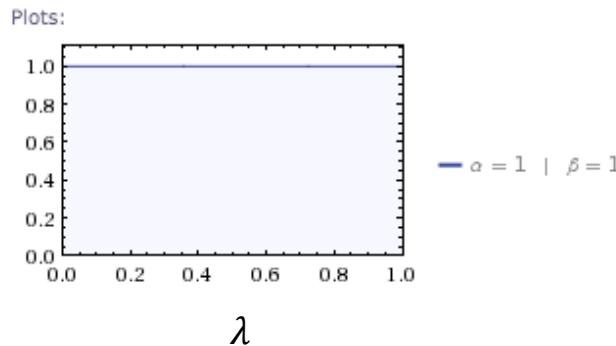
- Bernoulli distribution



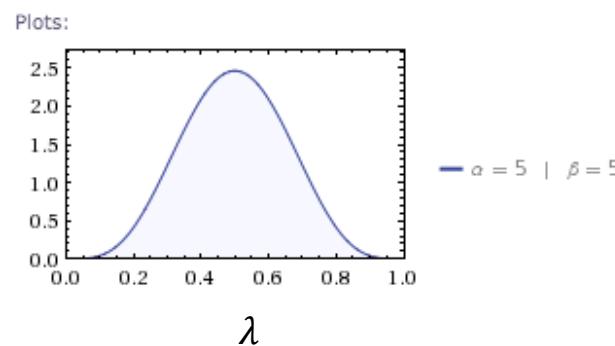
$$\prod_{i=1}^N P(x_n) = \prod_{i=1}^N \lambda^{x_n} (1 - \lambda)^{1-x_n}$$

Maximum a posteriori estimation (MAP): Bernoulli

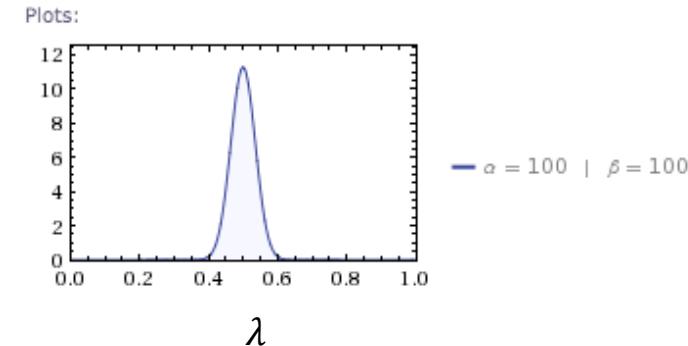
Beta(1,1) (λ)



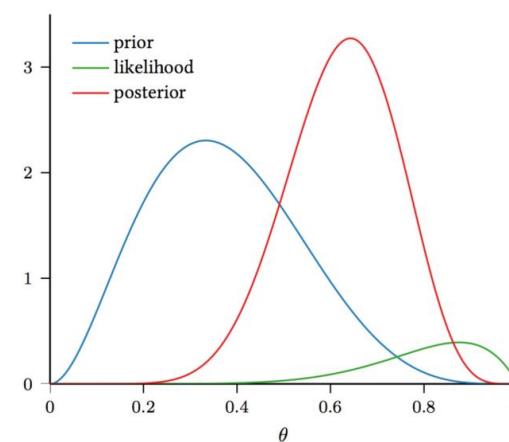
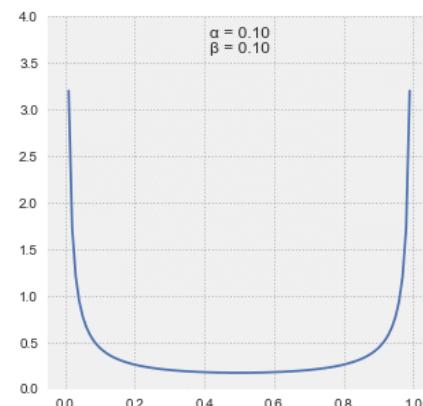
Beta(5,5) (λ)



Beta(100,100) (λ)



$$P(\lambda) = \frac{\lambda^{\alpha-1} (1-\lambda)^{\beta-1}}{Z(\alpha, \beta)}$$



$$P(\lambda|\mathcal{D}) \propto P(\mathcal{D}|\lambda)p(\lambda)$$

Maximum a posteriori estimation (MAP)

- $p(\mathcal{D}|\lambda)p(\lambda) \propto \lambda^{\# \text{ of heads}} \times (1 - \lambda)^{\# \text{ of tails}} \times \lambda^{\alpha-1} \times (1 - \lambda)^{\beta-1}$
 $= \lambda^{\# \text{ of heads} + \alpha - 1} \times (1 - \lambda)^{\# \text{ of tails} + \beta - 1}$
- What is $\underset{\lambda}{\operatorname{argmax}} P(\lambda|\mathcal{D}) = \underset{\lambda}{\operatorname{argmax}} \log p(\mathcal{D}|\lambda) + \log p(\lambda)$?
- $\circ \lambda = \frac{\# \text{ of heads} + \alpha - 1}{\# \text{ of samples} + \alpha + \beta - 2}$
- $\alpha - 1, \beta - 1$ (or α, β) are called pseudo counts
 - When N is large, MAP is similar to MLE

Comparison: MLE, MAP, and Bayesian

- Given “NH heads and NT tails”
 - MLE: $\underset{\lambda}{\operatorname{argmax}} P(\mathcal{D}|\lambda) = \frac{NH}{NH+NT}$
 - MAP: $\underset{\lambda}{\operatorname{argmax}} P(\lambda|\mathcal{D}) = \underset{\lambda}{\operatorname{argmax}} P(\mathcal{D}|\lambda)P(\lambda) = \frac{NH+\alpha-1}{(NH+\alpha-1)+(NT+\beta-1)}$, if $P(\lambda)$ is a Beta Dist.
 - Bayesian estimates: just **keep $P(\lambda|\mathcal{D})$** , with λ still a random variable
 - This keeps more information about λ (not just the mode)
 - MLE and MAP are point estimates: choosing one λ ; Bayesian: keep a distribution
- Usages after seeing \mathcal{D} : future prediction X = “head”
 - $P(X = \text{"head"}; \lambda) = \lambda$ λ is estimated by MLE or MAP
 - $P(X = \text{"head"}|\mathcal{D}) = \int P(X = \text{"head"}|\lambda) P(\lambda|\mathcal{D}) d\lambda = \frac{NH+\alpha}{(NH+\alpha)+(NT+\beta)}$; **“Bayesian inference”**

More about Bayesian

- $P(X = \text{"head"} | \mathcal{D}) = \int P(X = \text{"head"}, \lambda | D) d\lambda =$
- $\int P(X = \text{"head"} | \lambda, D) P(\lambda | \mathcal{D}) d\lambda = \int P(X = \text{"head"} | \lambda) P(\lambda | \mathcal{D}) d\lambda =$
- $\int \lambda P(\lambda | \mathcal{D}) d\lambda = E[\lambda | \mathcal{D}] = \frac{NH + \alpha}{(NH + \alpha) + (NT + \beta)}$
- In general (for other distributions), the integration is intractable, and we must approximate it by Monte Carlo methods, variational methods, etc.
 - $\int P(X=x | \lambda) P(\lambda | \mathcal{D}) d\lambda \approx \frac{1}{M} \sum_m P(X=x | \lambda_m)$, where $\lambda_m \sim P(\lambda | \mathcal{D})$
- $P(\lambda | \mathcal{D}) = \frac{P(\lambda, \mathcal{D})}{P(\mathcal{D})}$ is usually intractable, due to $P(\mathcal{D}) = \int P(\mathcal{D} | \lambda) P(\lambda) d\lambda$

Today

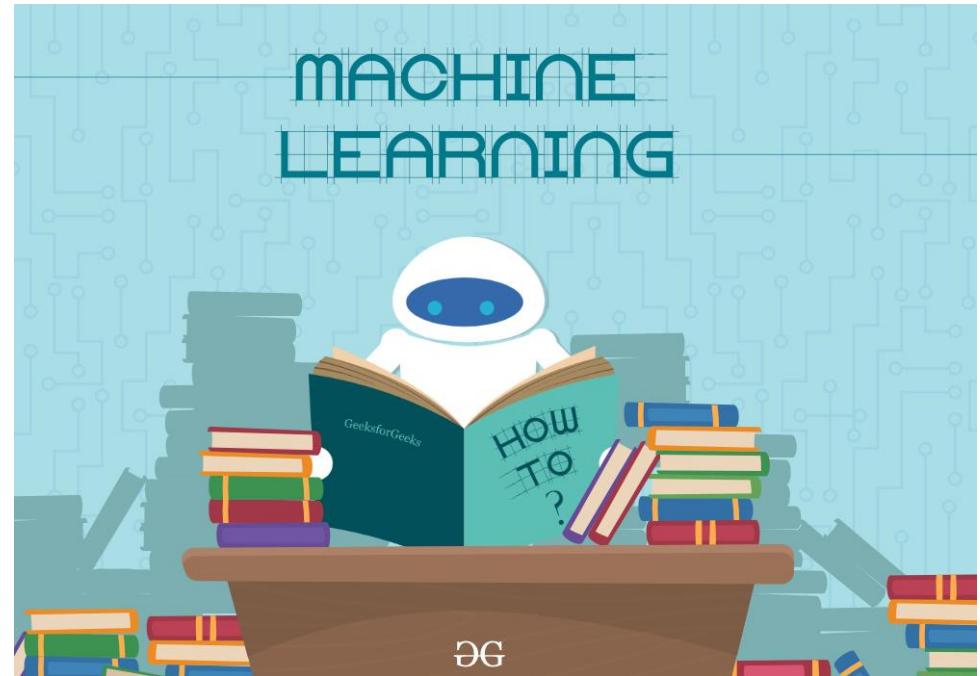
Review

Multi-dimensional Gaussian

Generative vs. discriminative learning

Linear regression revisit

Regularization, cross-validation



For multi-dimensional data: continuous

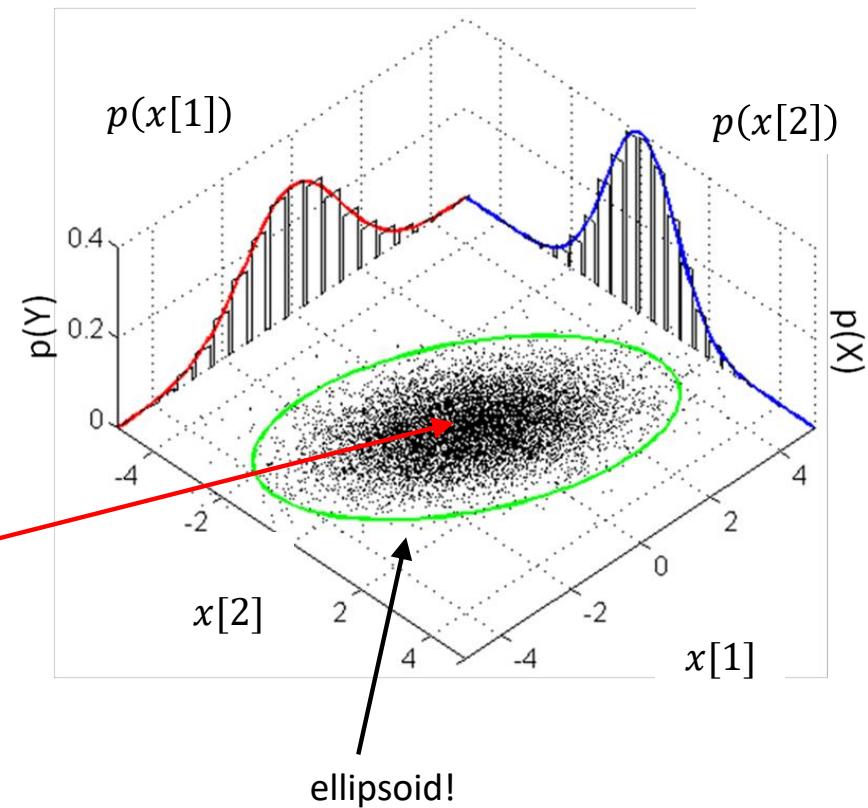
- Observe $\{x_i \in \mathbb{R}^D\}_{i=1}^N$
- Fit a multi-dimensional Gaussian (by MLE):

$$p(x) = \frac{1}{(2\pi)^{D/2} \det(\boldsymbol{\Sigma})^{1/2}} e^{-\frac{1}{2}(x-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (x-\boldsymbol{\mu})}$$

$$\underset{\boldsymbol{\mu}, \boldsymbol{\Sigma}}{\operatorname{argmax}} p(x_1, \dots, x_N; \quad \boldsymbol{\mu} \in \mathbb{R}^D, \boldsymbol{\Sigma} \in \mathbb{R}^{D \times D})$$

$$\bullet \boldsymbol{\mu} = \frac{\sum_{i=1}^N x_i}{N}$$

$$\bullet \boldsymbol{\Sigma} = \text{covariance matrix} = \frac{1}{N} \sum_{i=1}^N (x_i - \boldsymbol{\mu})(x_i - \boldsymbol{\mu})^T$$



Background: covariance matrix

$$\begin{aligned}\Sigma &= \begin{bmatrix} \text{cov}_{1,1} & \dots & \text{cov}_{1,D} \\ \vdots & \ddots & \vdots \\ \text{cov}_{D,1} & \dots & \text{cov}_{D,D} \end{bmatrix} = \begin{bmatrix} E[(x[1] - \mu[1])(x[1] - \mu[1])] & \dots & E[(x[1] - \mu[1])(x[D] - \mu[D])] \\ \vdots & \ddots & \vdots \\ E[(x[D] - \mu[D])(x[1] - \mu[1])] & \dots & E[(x[D] - \mu[D])(x[D] - \mu[D])] \end{bmatrix} \\ &= \begin{bmatrix} \frac{1}{N} \sum_{n=1}^N [(x_n[1] - \mu[1])(x_n[1] - \mu[1])] & \dots & \frac{1}{N} \sum_{n=1}^N [(x_n[1] - \mu[1])(x_n[D] - \mu[D])] \\ \vdots & \ddots & \vdots \\ \frac{1}{N} \sum_{n=1}^N [(x_n[D] - \mu[D])(x_n[1] - \mu[1])] & \dots & \frac{1}{N} \sum_{n=1}^N [(x_n[D] - \mu[D])(x_n[D] - \mu[D])] \end{bmatrix} \\ &= \frac{1}{N} \sum_{n=1}^N \begin{bmatrix} (x_n[1] - \mu[1])(x_n[1] - \mu[1]) & \dots & (x_n[1] - \mu[1])(x_n[D] - \mu[D]) \\ \vdots & \ddots & \vdots \\ (x_n[D] - \mu[D])(x_n[1] - \mu[1]) & \dots & (x_n[D] - \mu[D])(x_n[D] - \mu[D]) \end{bmatrix} \\ &= \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \boldsymbol{\mu})(\mathbf{x}_n - \boldsymbol{\mu})^T \\ &= \frac{1}{N} \bar{\mathbf{X}} \bar{\mathbf{X}}^T, \text{ where } \bar{\mathbf{X}} = [\mathbf{x}_1 - \boldsymbol{\mu}, \quad \mathbf{x}_2 - \boldsymbol{\mu}, \quad \mathbf{x}_3 - \boldsymbol{\mu}, \quad \dots, \quad \mathbf{x}_N - \boldsymbol{\mu}]\end{aligned}$$

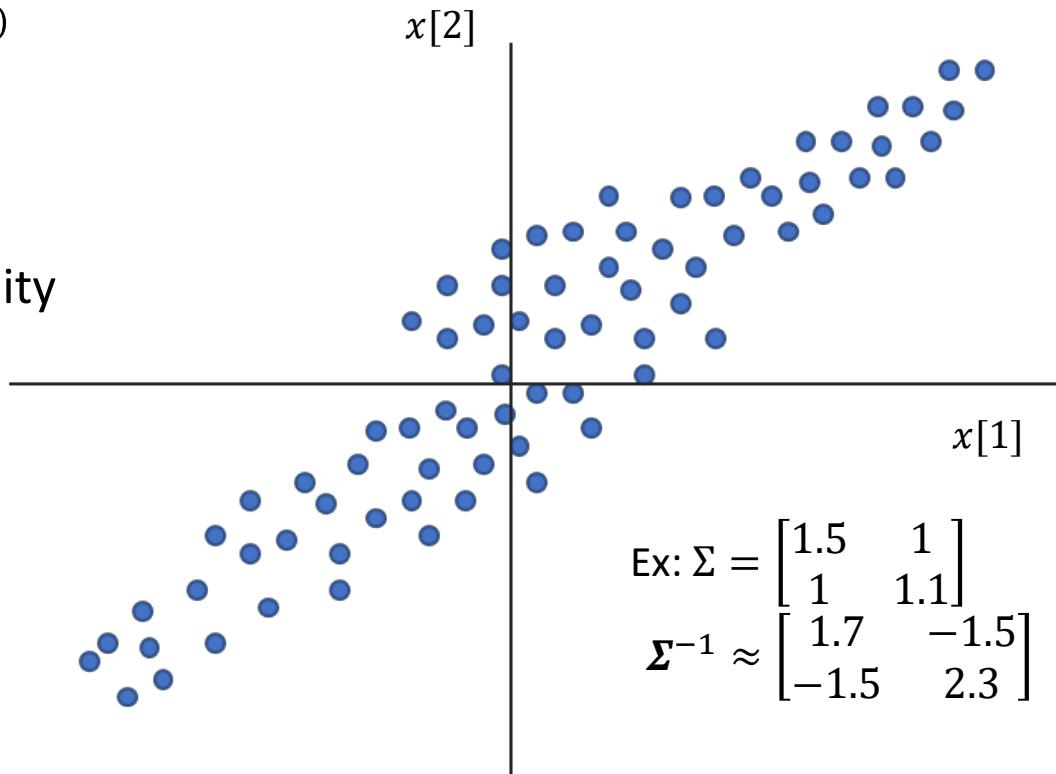
One-dimensional vs. multi-dimensional

- One-dimensional: $p(X = x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/2\sigma^2}$
- Multi-dimensional: $p(x) = \frac{1}{(2\pi)^{D/2} \det(\Sigma)^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$
- When $\Sigma = \sigma^2 I$, $p(x) = \frac{1}{(2\pi)^{D/2} \sigma^D} e^{-(x-\mu)^T (x-\mu)/2\sigma^2} = \frac{1}{(2\pi)^{D/2} \sigma^D} e^{-\|x-\mu\|_2^2/2\sigma^2}$

For multi-dimensional data: continuous

$$p(x) = \frac{1}{(2\pi)^{D/2} \det(\Sigma)^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

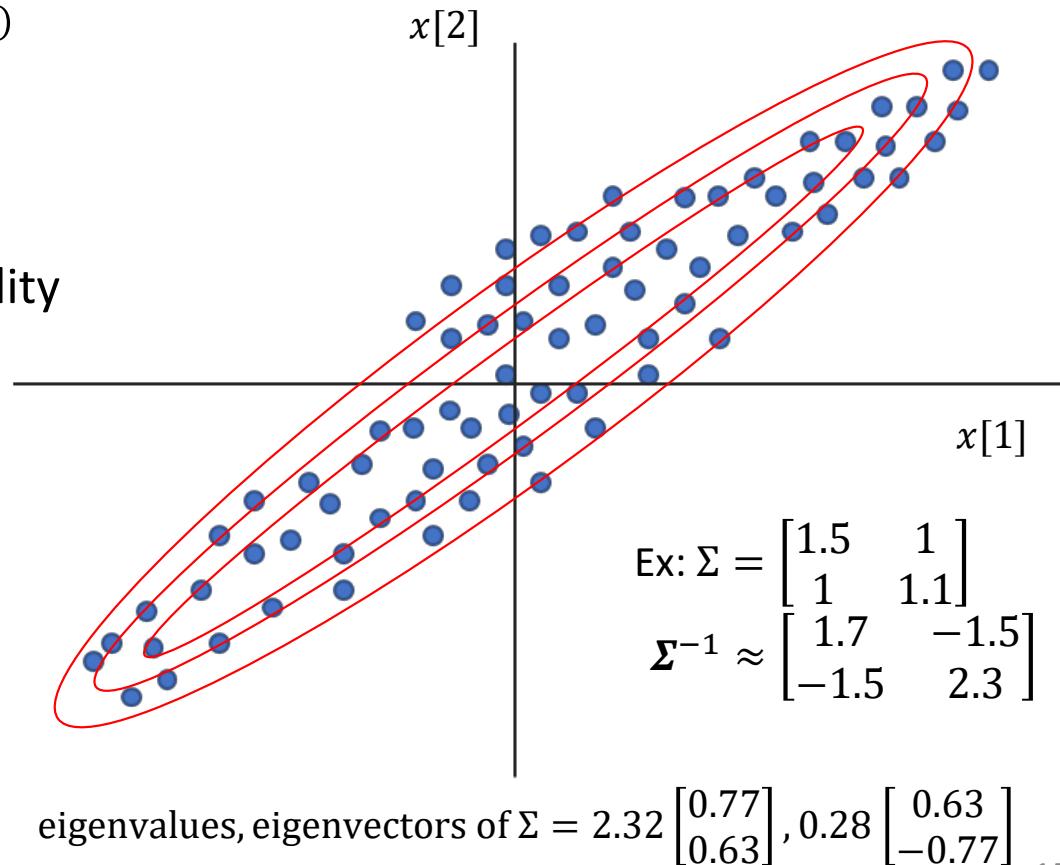
- Properties:
 - Faraway from the mean, lower probability



For multi-dimensional data: continuous

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{D/2} \det(\boldsymbol{\Sigma})^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu})}$$

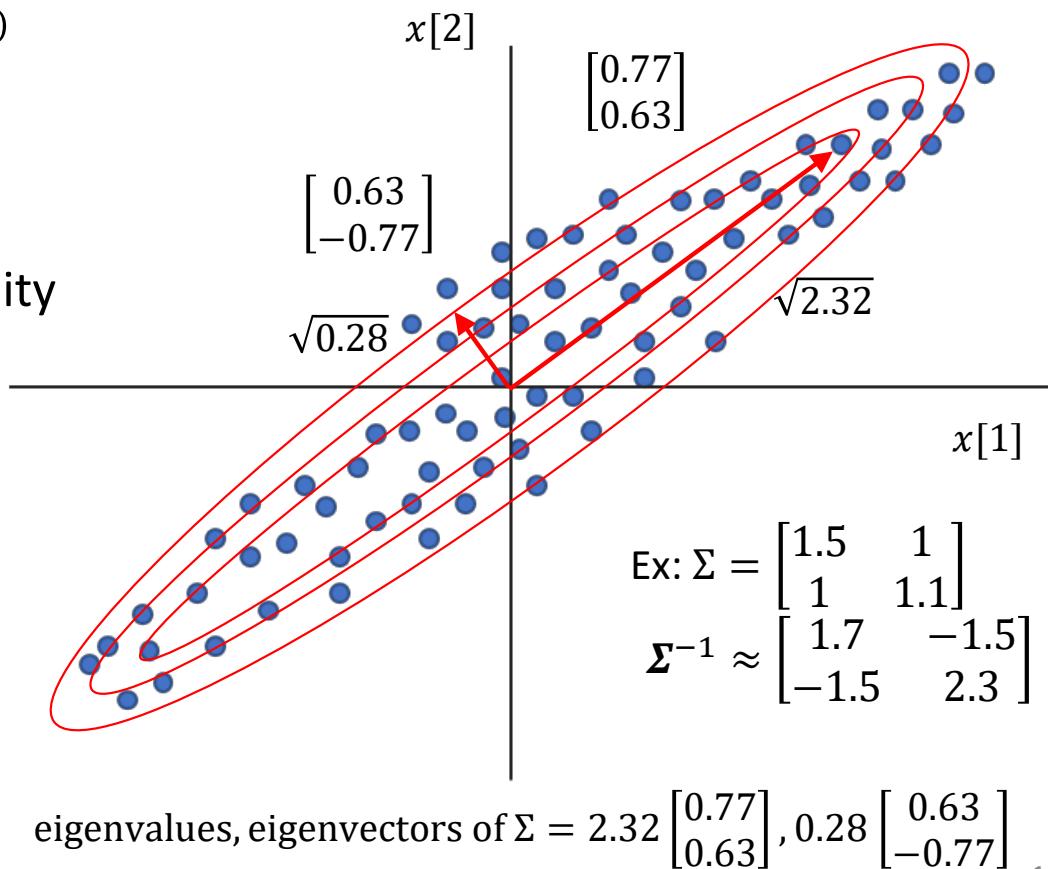
- Properties:
 - Faraway from the mean, lower probability
- Think
 - What is $\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$?



For multi-dimensional data: continuous

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{D/2} \det(\boldsymbol{\Sigma})^{1/2}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu})}$$

- Properties:
 - Faraway from the mean, lower probability
- Think
 - What is $\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$?



Today

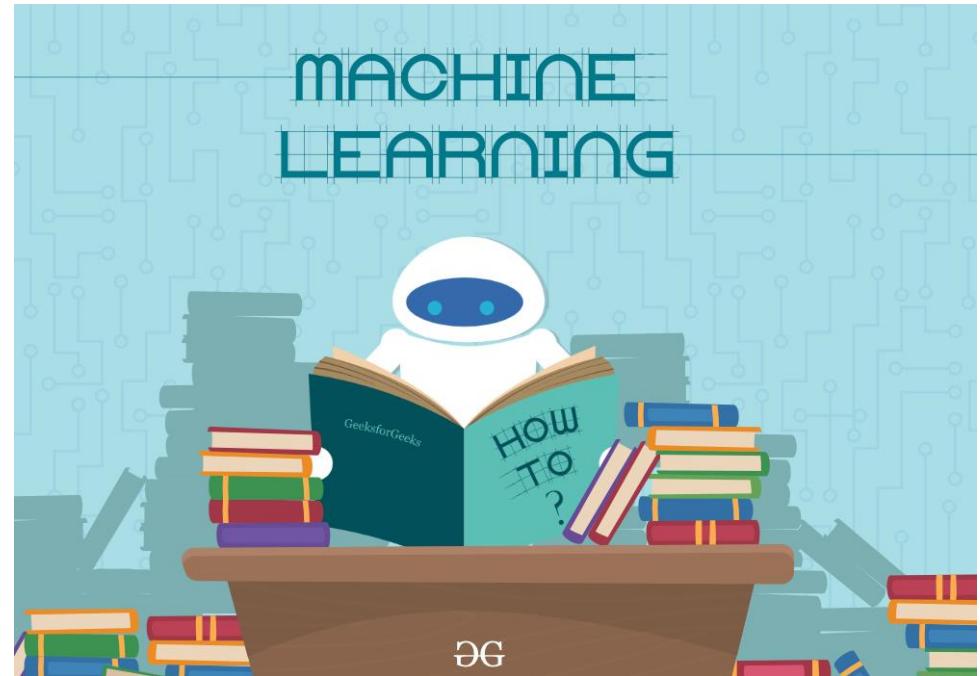
Review

Multi-dimensional Gaussian

Generative vs. discriminative learning

Linear regression revisit

Regularization, cross-validation



Generative vs. Discriminative learning

- Training data: $D_{tr} = \{(x_n, y_n)\}_{n=1}^N$
- Goal: fit $p(X, Y)$ to D_{tr}
 - $p(Y)p(X|Y)$: generative learning (with $p(X|Y)$ you can generate/sample new data)
 - $p(X)p(Y|X)$: discriminative learning
 - Usually not estimating $p(X)$
- Prediction (in the MLE or MAP scenario, mostly we do)
 - Given a test example x
 - $\hat{y} = \max_c p(Y = c|x)$; the Bayes Optimal Classifier
 - In this sense, discriminative learning seems to be more straightforward (or specialized)

Generative vs. Discriminative learning

- Generative learning (MLE case)
 - $p(x, y; \theta) = p(x|y; \theta)p(y; \theta)$
 - $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} p(D_{tr}; \theta) = \underset{\theta}{\operatorname{argmax}} \prod_n p(x_n, y_n ; \theta) = \underset{\theta}{\operatorname{argmax}} \prod_n p(x_n|y_n ; \theta)p(y_n; \theta)$
 - Need to make two modeling assumptions (conditional data distribution, prior)
 - Modeling $p(x_n|y_n ; \theta)$ is challenging
 - If the assumption is correct, then fewer training data are needed

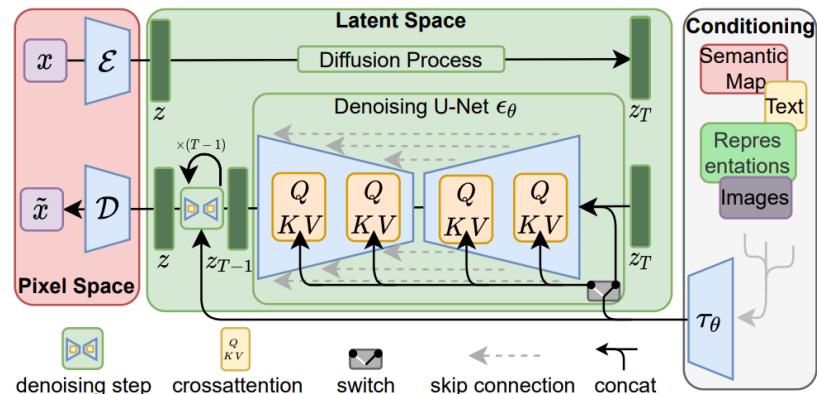
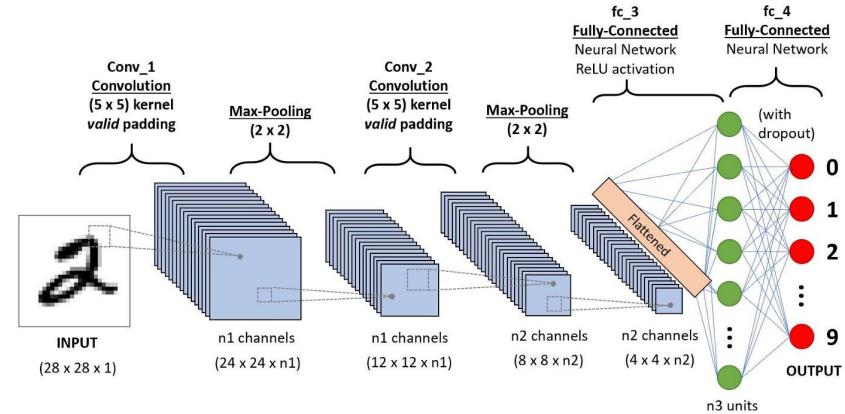


Generative vs. Discriminative learning

- Generative learning (MLE case)
 - $p(x, y; \theta) = p(x|y; \theta)p(y; \theta)$
 - $\widehat{\theta} = \underset{\theta}{\operatorname{argmax}} p(D_{tr}; \theta) = \underset{\theta}{\operatorname{argmax}} \prod_n p(x_n, y_n; \theta) = \underset{\theta}{\operatorname{argmax}} \prod_n p(x_n|y_n; \theta)p(y_n; \theta)$
 - Need to make two modeling assumptions (conditional data distribution, prior)
 - Modeling $p(x_n|y_n; \theta)$ is challenging
 - If the assumption is correct, then fewer training data are needed
- Discriminative learning (MLE case)
 - $p(x, y; \theta) = p(y|x; \theta)p(x)$
 - $\widehat{\theta} = \underset{\theta}{\operatorname{argmax}} p(D_{tr}; \theta) = \underset{\theta}{\operatorname{argmax}} \prod_n p(y_n|x_n; \theta)$
 - Need to make one modeling assumptions (conditional class distribution)

Don't get confused

- Discriminative vs. generative learning:
 - Difference in **modeling**
- MAP vs. MLE:
 - Difference in **the objective function**



Today

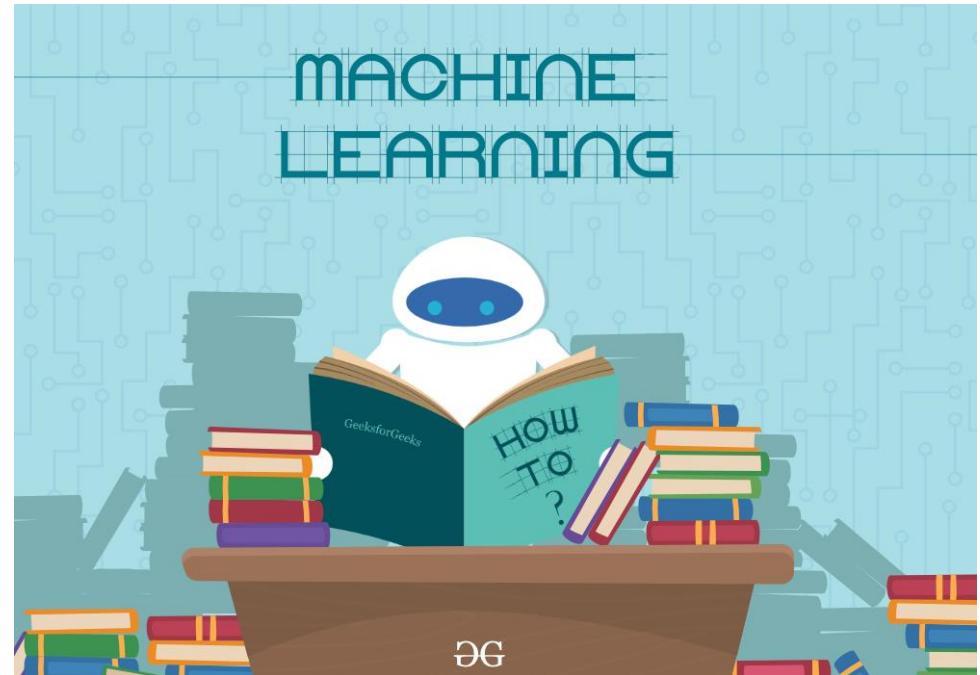
Review

Multi-dimensional Gaussian

Generative vs. discriminative learning

Linear regression revisit

Regularization, cross-validation

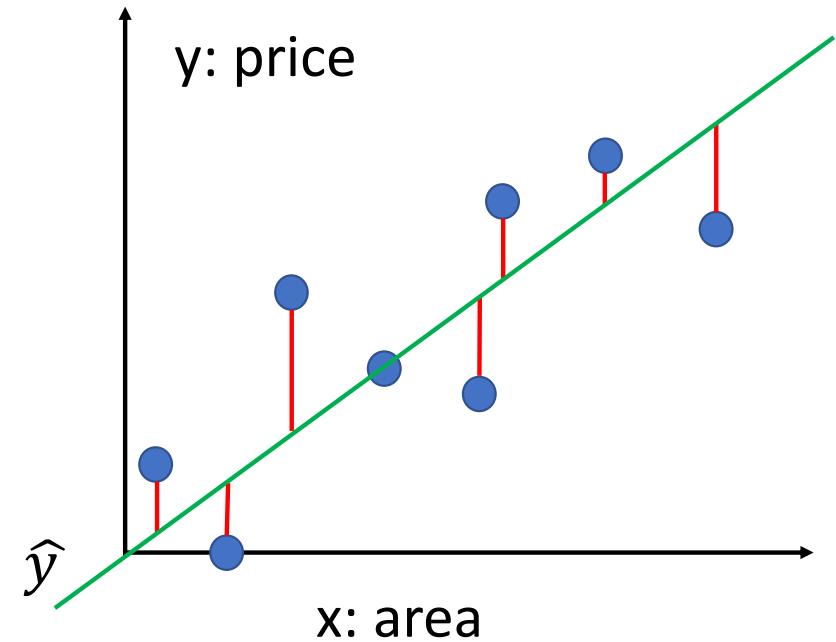


Linear regression revisit (bias is absorbed)

- Given $D_{tr} = \{(\mathbf{x}_n, y_n \in \mathbb{R})\}_{n=1}^N$
- Goal: Fit a probabilistic model that can be used to predict label given \mathbf{x}
- Let us consider $p(X, Y) = p(X)p(Y|X; \boldsymbol{\theta})$ or $p(X, Y) = p(X)p(Y|X, \boldsymbol{\theta})$
 - We assume that $p(X)$ is unknown but fixed (so no need to learn it)
 - We will only estimate the parameters $\hat{\boldsymbol{\theta}}$ for $p(Y|X; \boldsymbol{\theta})$
 - $\underset{y}{\operatorname{argmax}} p(Y = y | X = \mathbf{x}; \hat{\boldsymbol{\theta}})$ for prediction, for example
- $p(D_{tr}; \boldsymbol{\theta}) = p(\{(\mathbf{x}_n, y_n)\}_{n=1}^N; \boldsymbol{\theta}) = \prod_n p(\mathbf{x}_n) \prod_n p(y_n | \mathbf{x}_n; \boldsymbol{\theta})$
 - That is, $p(D_{tr}; \boldsymbol{\theta}) \propto \prod_n p(y_n | \mathbf{x}_n; \boldsymbol{\theta})$

Linear regression revisit (bias is absorbed)

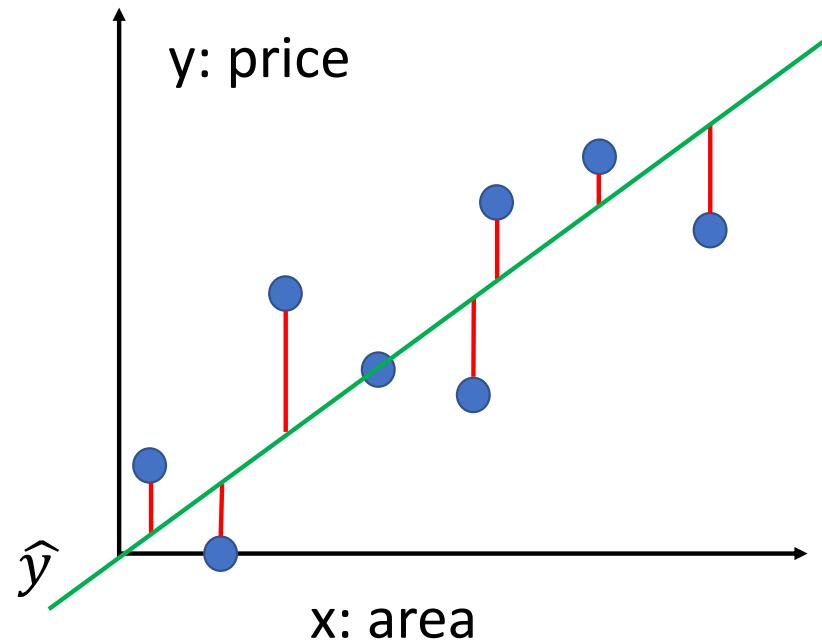
- Given $D_{tr} = \{(x_n, y_n) \in \mathbb{R}^2\}_{n=1}^N$
- Goal: We want to fit $p(Y|X; \theta)$
- What probabilistic model should we use?
- We use: $p(y|x; w) \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-w^T x)^2/2\sigma^2}$
 - $Y = y$ is a normally distributed around $w^T x$



Linear regression revisit (bias is absorbed)

- Assume that $y = \mathbf{w}^T \mathbf{x} + \varepsilon$, where $\varepsilon \sim \mathcal{N}(0, \sigma^2)$
 - $\mathcal{N}(\mu, \sigma^2)$ has a PDF $p(\varepsilon) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(\varepsilon-0)^2/2\sigma^2}$
 - y is drawn from a line $\mathbf{w}^T \mathbf{x}$ plus some Gaussian noise

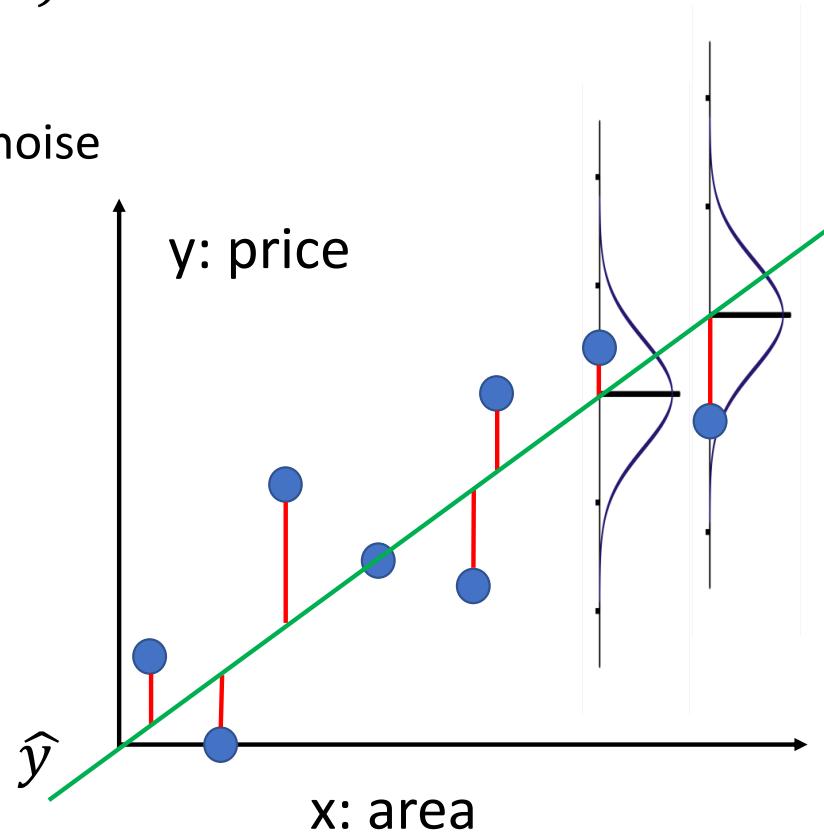
- $y|x \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}, \sigma^2)$
- $p(y|x; \mathbf{w}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-\mathbf{w}^T \mathbf{x})^2/2\sigma^2}$



Linear regression revisit (bias is absorbed)

- Assume that $y = \mathbf{w}^T \mathbf{x} + \varepsilon$, where $\varepsilon \sim \mathcal{N}(0, \sigma^2)$
 - $\mathcal{N}(\mu, \sigma^2)$ has a PDF $p(\varepsilon) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(\varepsilon-0)^2/2\sigma^2}$
 - y is drawn from a line $\mathbf{w}^T \mathbf{x}$ plus some Gaussian noise

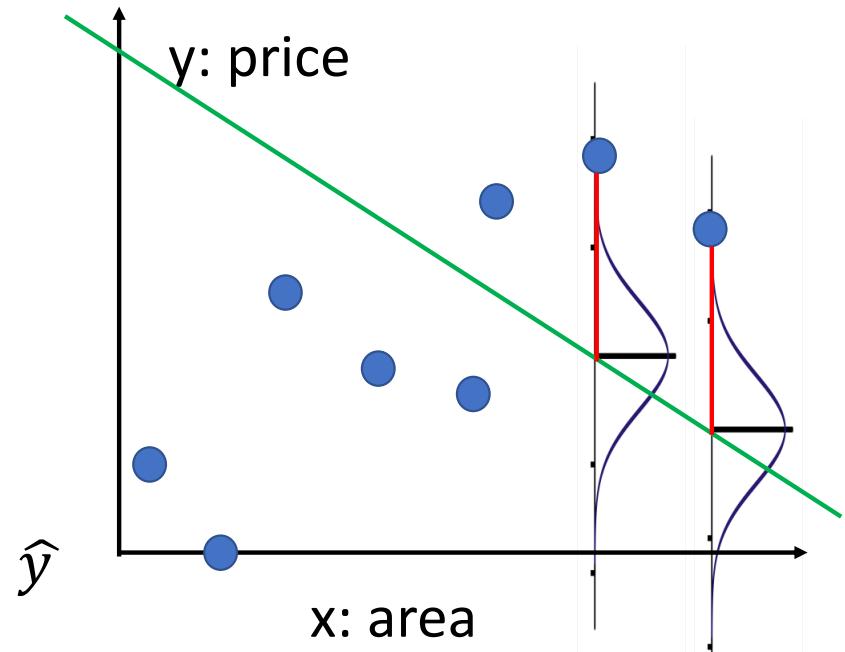
- $y|x \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}, \sigma^2)$
- $p(y|x; \mathbf{w}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-\mathbf{w}^T \mathbf{x})^2/2\sigma^2}$



Linear regression revisit (bias is absorbed)

- Assume that $y = \mathbf{w}^T \mathbf{x} + \varepsilon$, where $\varepsilon \sim \mathcal{N}(0, \sigma^2)$
 - $\mathcal{N}(\mu, \sigma^2)$ has a PDF $p(\varepsilon) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(\varepsilon-0)^2/2\sigma^2}$
 - y is drawn from a line $\mathbf{w}^T \mathbf{x}$ plus some Gaussian noise

- $y|x \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}, \sigma^2)$
- $p(y|x; \mathbf{w}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y-\mathbf{w}^T \mathbf{x})^2/2\sigma^2}$



Linear regression revisit (bias is absorbed)

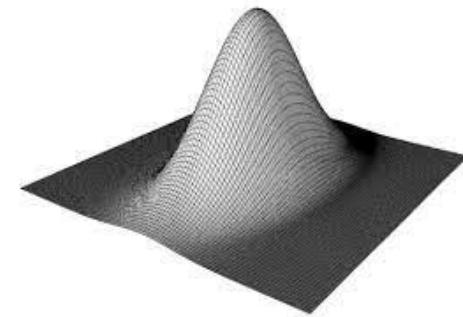
- Setup
 - $D_{tr} = \{(\mathbf{x}_n \in \mathbb{R}^D, y_n \in \mathbb{R})\}_{n=1}^N$
 - A fixed unknown $p(\mathbf{X} = \mathbf{x})$
- $p(D_{tr}; \mathbf{w}) = \prod_{n=1}^N p(\mathbf{x}_n, y_n; \mathbf{w}) = \prod_{n=1}^N p(\mathbf{x}_n)p(y_n | \mathbf{x}_n; \mathbf{w}) \propto \prod_{n=1}^N p(y_n | \mathbf{x}_n; \mathbf{w})$
- $\log \prod_{n=1}^N p(y_n | \mathbf{x}_n; \mathbf{w}) = \sum_n \log p(y_n | \mathbf{x}_n; \mathbf{w}) = \sum_n \log \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y_n - \mathbf{w}^T \mathbf{x}_n)^2 / 2\sigma^2} = \sum_n \log \frac{1}{\sqrt{2\pi\sigma^2}} + \log e^{-(y_n - \mathbf{w}^T \mathbf{x}_n)^2 / 2\sigma^2}$

Linear regression revisit (bias is absorbed)

- $\sum_n \log \frac{1}{\sqrt{2\pi\sigma^2}} + \log e^{-(y_n - \mathbf{w}^T \mathbf{x}_n)^2/2\sigma^2} \propto \sum_n -(y_n - \mathbf{w}^T \mathbf{x}_n)^2/2\sigma^2 \propto \sum_n -(y_n - \mathbf{w}^T \mathbf{x}_n)^2$
- MLE: $\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmax}} \sum_n -(y_n - \mathbf{w}^T \mathbf{x}_n)^2 = \underset{\theta}{\operatorname{argmin}} \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n)^2$

Ridge regression revisit

- Assume that $p(\mathbf{w}) = \frac{1}{(2\pi)^{D/2} \det(\tau^2 \mathbf{I})^{1/2}} e^{-\frac{1}{2}(\mathbf{w}-\mathbf{0})^T (\tau^2 \mathbf{I})^{-1} (\mathbf{w}-\mathbf{0})}$
- $p(\mathbf{w}) = \frac{1}{(2\pi)^{D/2} \det(\tau^2 \mathbf{I})^{1/2}} e^{-\frac{\|\mathbf{w}\|_2^2}{2\tau^2}}$
- $\log p(D_{tr}|\mathbf{w})p(\mathbf{w}) \propto \sum_n -\left(y_n - \mathbf{w}^T \mathbf{x}_n\right)^2 + \log p(\mathbf{w})$
- $\underset{\mathbf{w}}{\operatorname{argmax}} \sum_n -\left(y_n - \mathbf{w}^T \mathbf{x}_n\right)^2 - \frac{1}{2} \frac{\|\mathbf{w}\|_2^2}{\tau^2} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_n \left(y_n - \mathbf{w}^T \mathbf{x}_n\right)^2 + \frac{1}{2} \frac{\|\mathbf{w}\|_2^2}{\tau^2}$
- Ridge regression includes an L2 regularizer (not dependent on data)



MLE, MAP, and Bayesian for supervised learning

- MLE: $\widehat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(D_{tr}; \boldsymbol{\theta})$, and predict by $p(y|x; \widehat{\boldsymbol{\theta}})$
- MAP: $\widehat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\boldsymbol{\theta}|D_{tr}) \propto p(D_{tr}|\boldsymbol{\theta})p(\boldsymbol{\theta})$, and predict by $p(y|x; \widehat{\boldsymbol{\theta}})$
- Difference between MAP and MLE is subtle
 - $\log p(\boldsymbol{\theta})$ can be seen as a regularizer to penalize largely deviating $\boldsymbol{\theta}$
- Bayesian: find $p(\boldsymbol{\theta}|D_{tr})$ and predict by $p(y|x, D_{tr}) = \int p(y|x, \boldsymbol{\theta}) P(\boldsymbol{\theta}|D_{tr}) d\boldsymbol{\theta}$

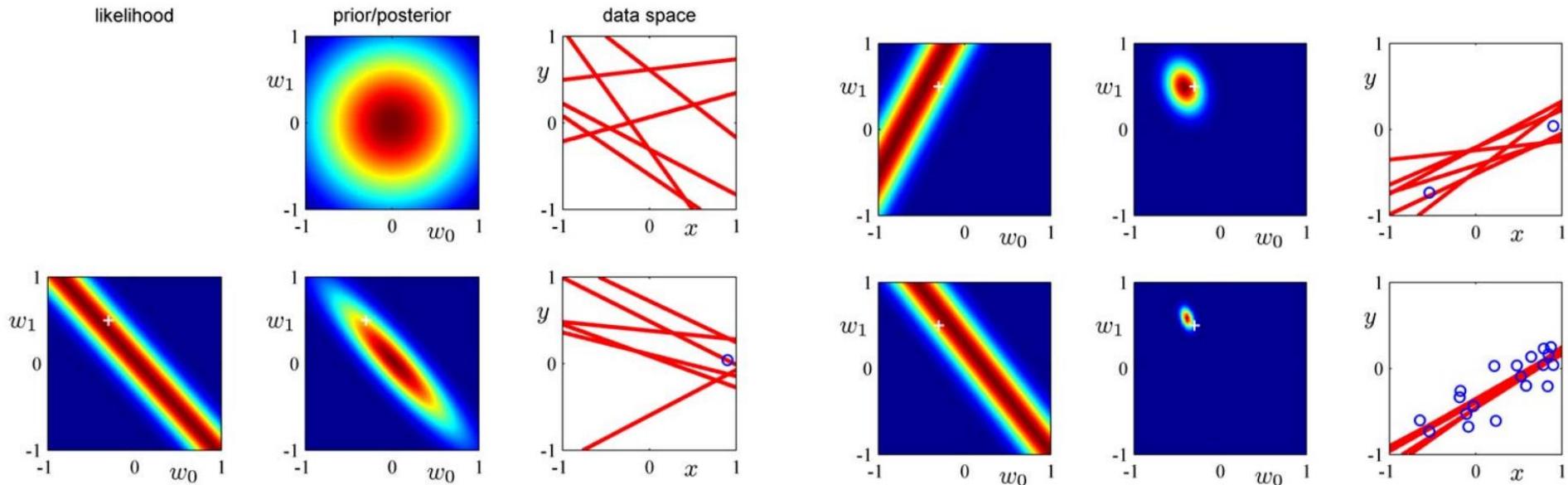
Questions?



THE OHIO STATE UNIVERSITY

* Bayesian (linear/nonlinear) regression *

- $p(y|x, D_{tr}) = \int p(y|x, w) p(w|D_{tr}) d\lambda$
- Let $p(y|x, w) = \mathcal{N}(w[0] + w[1]x, \sigma^2)$, $p(w) = \mathcal{N}(\mathbf{0}, I)$

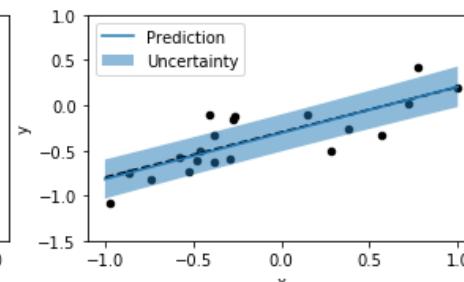
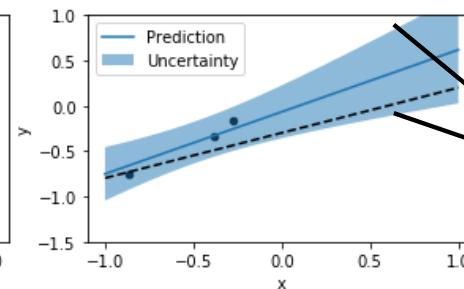
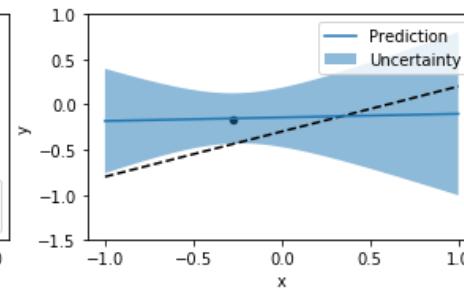
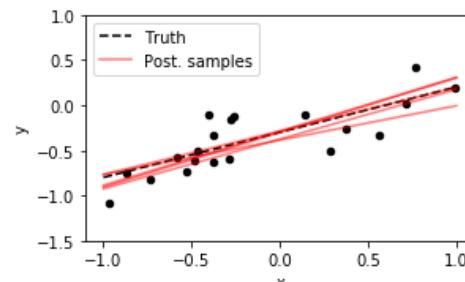
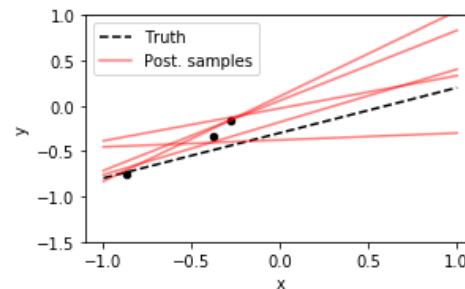
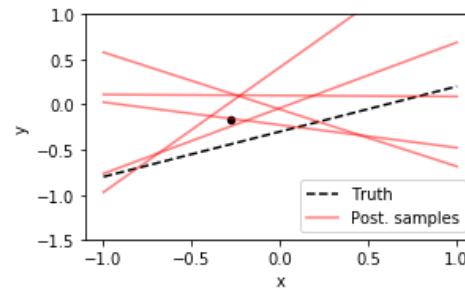
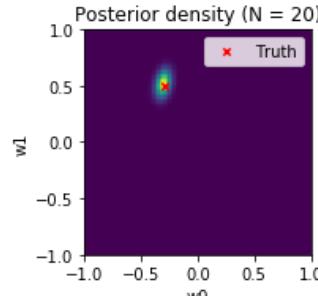
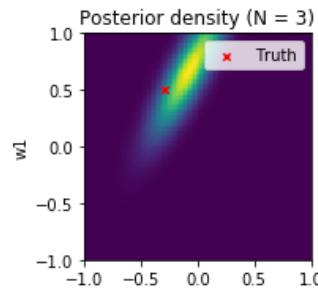
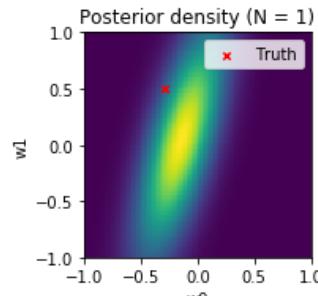


[Figure by Bishop, PRML, 2006]

* Bayesian (linear/nonlinear) regression *

$$p(\mathbf{w}|D_{tr})$$

$$p(y|\mathbf{x}, \mathbf{w}) = \mathcal{N}(\mathbf{w}^T \mathbf{x}, \sigma^2)$$



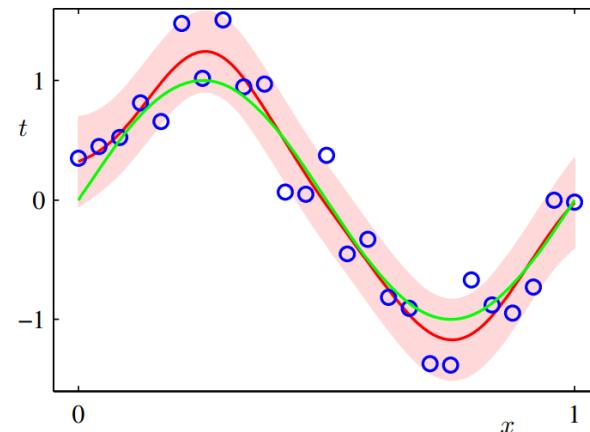
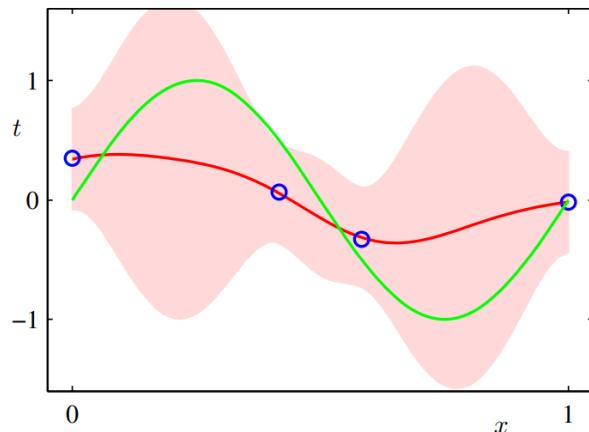
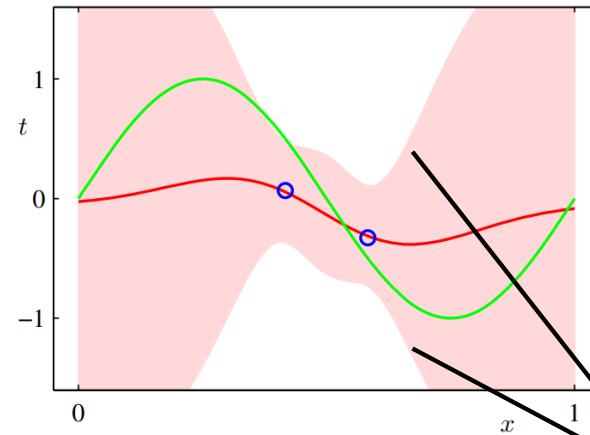
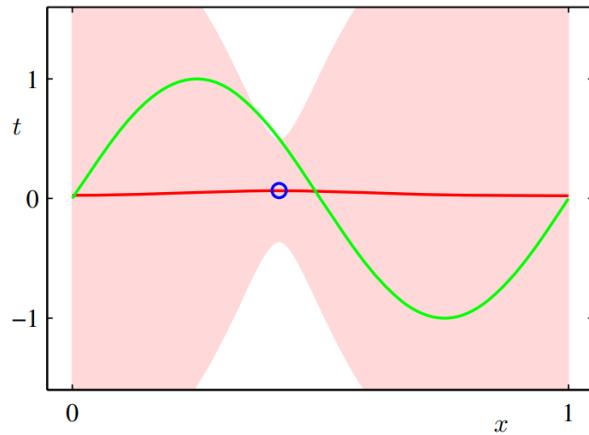
$$p(y|\mathbf{x}, D_{tr}) =$$

$$\int p(y|\mathbf{x}, \mathbf{w}) p(\mathbf{w}|D_{tr}) d\lambda$$

Sample-dependent
prediction variance

* Bayesian (linear/nonlinear) regression *

With feature transform



$$p(y|x, D_{tr}) = \int p(y|x, w) p(w|D_{tr}) d\lambda$$

Sample-dependent prediction variance

Today

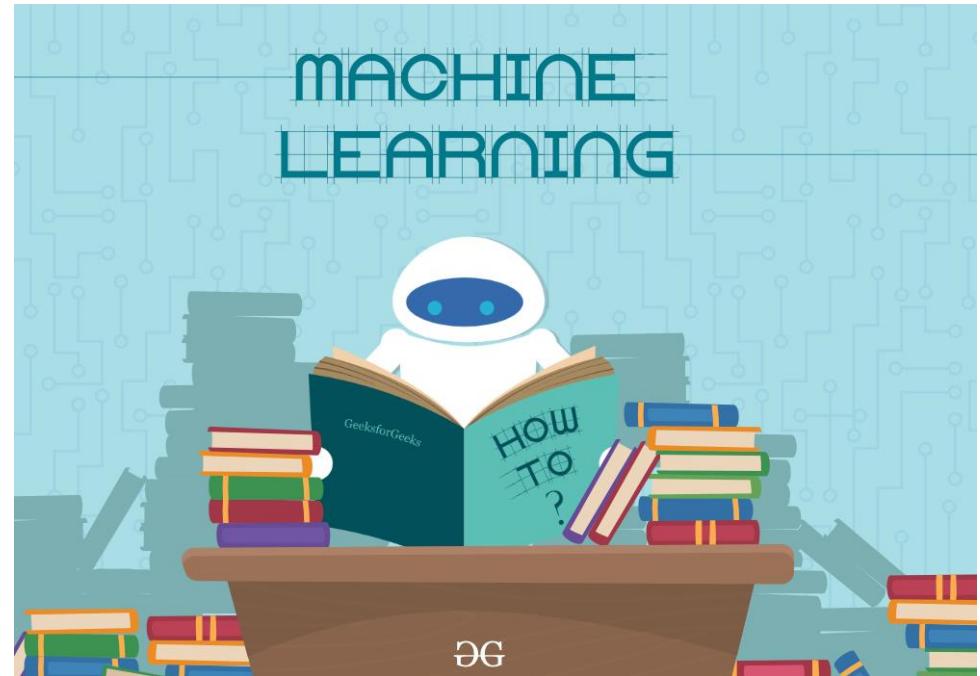
Review

Multi-dimensional Gaussian

Generative vs. discriminative learning

Linear regression revisit

Regularization, cross-validation



Regularization and hyper-parameters

- Ridge regression:

- $\underset{\mathbf{w}}{\operatorname{argmin}} \quad \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \lambda \|\mathbf{w}\|_2^2$

Regularization and hyper-parameters

- Ridge regression:

- $\underset{\mathbf{w}}{\operatorname{argmin}} \quad \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \lambda \|\mathbf{w}\|_2^2$

- The object function contains a training loss term and a regularization (e.g., belief) term

Regularization and hyper-parameters

- Ridge regression:

- $\underset{\mathbf{w}}{\operatorname{argmin}} \quad \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \lambda \|\mathbf{w}\|_2^2$

- The **object function** contains a **training loss term** and a **regularization (e.g., belief) term**
- Model parameters (e.g., \mathbf{w}) is learned from minimizing the **training objective**

Regularization and hyper-parameters

- Ridge regression:

$$\underset{\mathbf{w}}{\operatorname{argmin}} \quad \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \lambda \|\mathbf{w}\|_2^2$$

- The **object function** contains a **training loss term** and a **regularization (e.g., belief) term**
- Model parameters (e.g., \mathbf{w}) is learned from minimizing the **training objective**
- How about hyper-parameters (e.g., λ) of the algorithm or hypothesis class?
 - Can you determine λ by minimizing the training objective?

Regularization and hyper-parameters

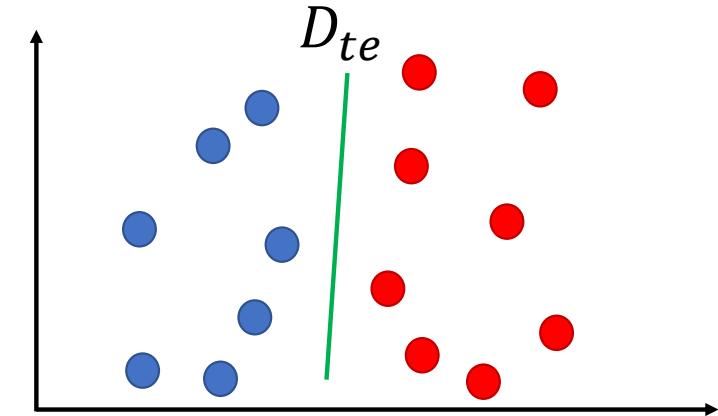
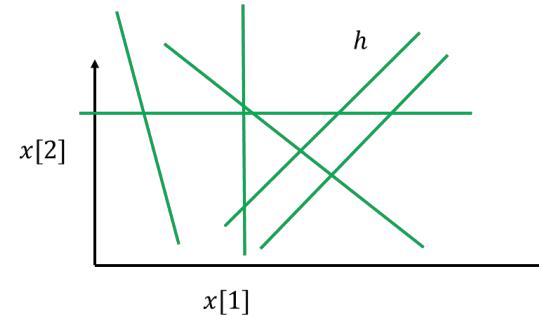
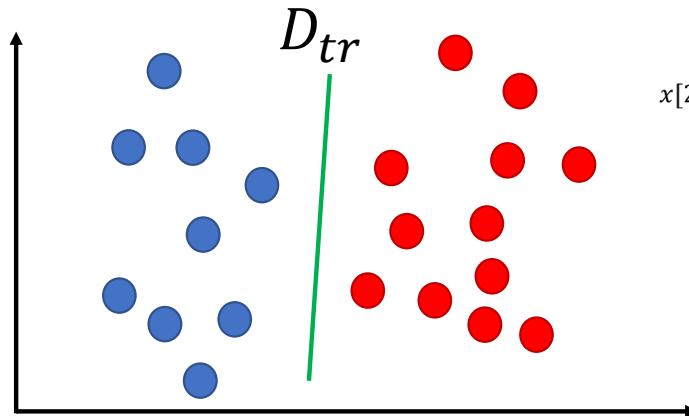
- Ridge regression:

$$\underset{\mathbf{w}}{\operatorname{argmin}} \quad \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \lambda \|\mathbf{w}\|_2^2$$

- The **object function** contains a **training loss term** and a **regularization (e.g., belief) term**
- Model parameters (e.g., \mathbf{w}) is learned from minimizing the **training objective**
- How about hyper-parameters (e.g., λ) of the algorithm or hypothesis class?
 - Can you determine λ by minimizing the training objective?
- Caution: Don't be confused by notations; λ is not always used as parameters.

Machine learning

- A good learning algorithm + hypothesis class
 - Contain a hypothesis that can achieve low **training loss** and **test (validation) loss**
 - The algorithm can return such a hypothesis
 - We may view ridge regression as regularized loss (risk) minimization + linear regressors



- In practice, we choose the learning algorithm + hypothesis class by validation

(Cross) validation

- Goal: choose the hyper-parameters (ID of an algorithm + hypothesis class)
 - Ex:# of neighbors and distance in KNN, λ in ridge regression, η in iterative optimization

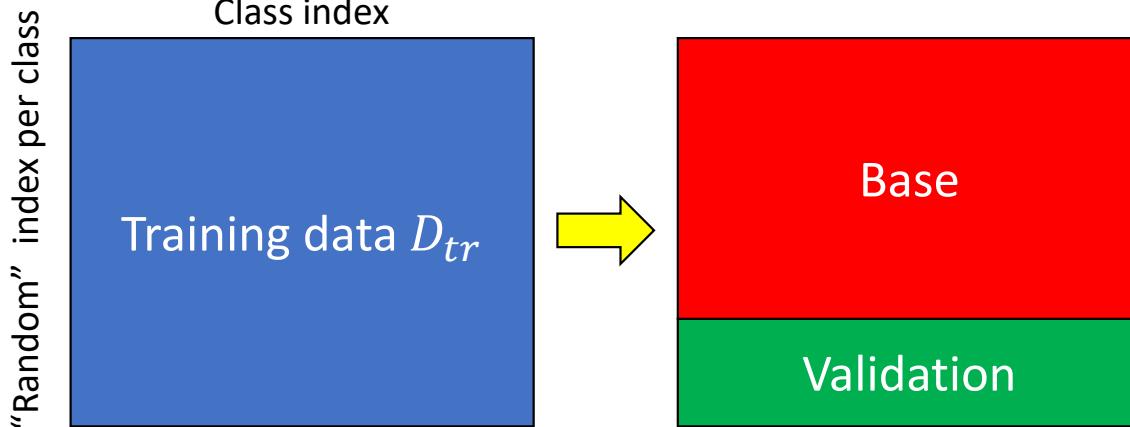


(Cross) validation

- Goal: choose the hyper-parameters
 - Ex:# of neighbors and distance in KNN, λ in ridge regression, η in iterative optimization

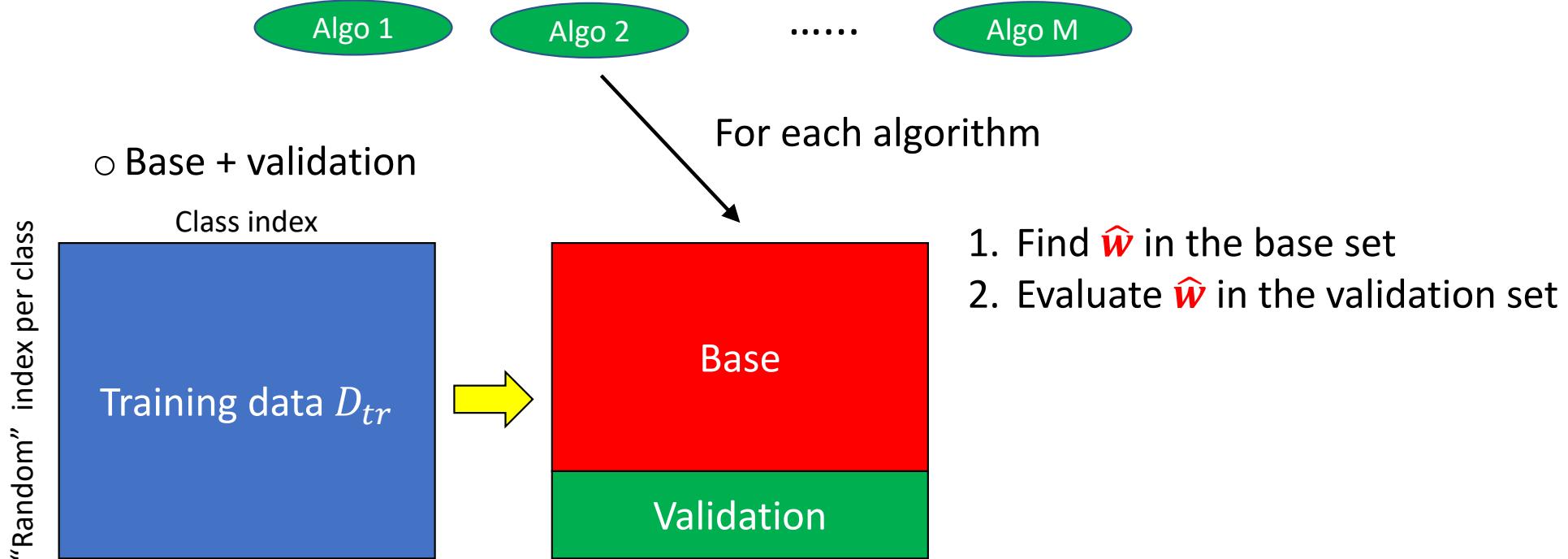


- Base + validation



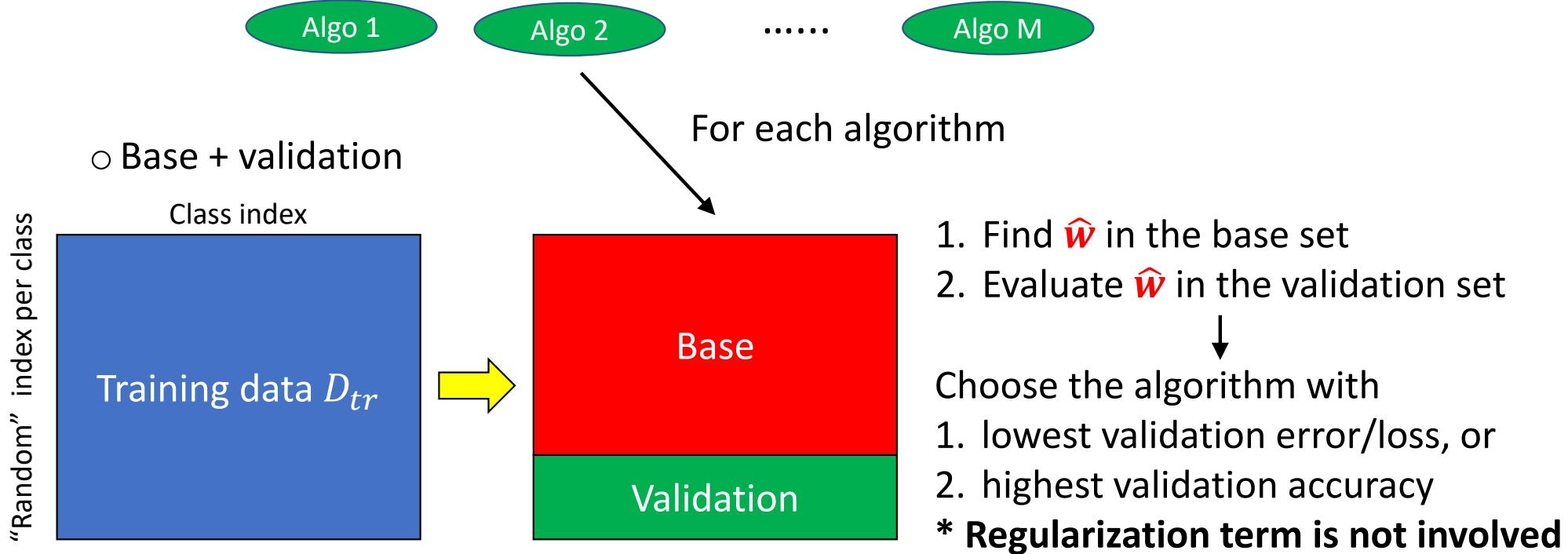
(Cross) validation

- Goal: choose the hyper-parameters
 - Ex:# of neighbors and distance in KNN, λ in ridge regression, η in iterative optimization



(Cross) validation

- Goal: choose the hyper-parameters
 - Ex:# of neighbors and distance in KNN, λ in ridge regression, η in iterative optimization

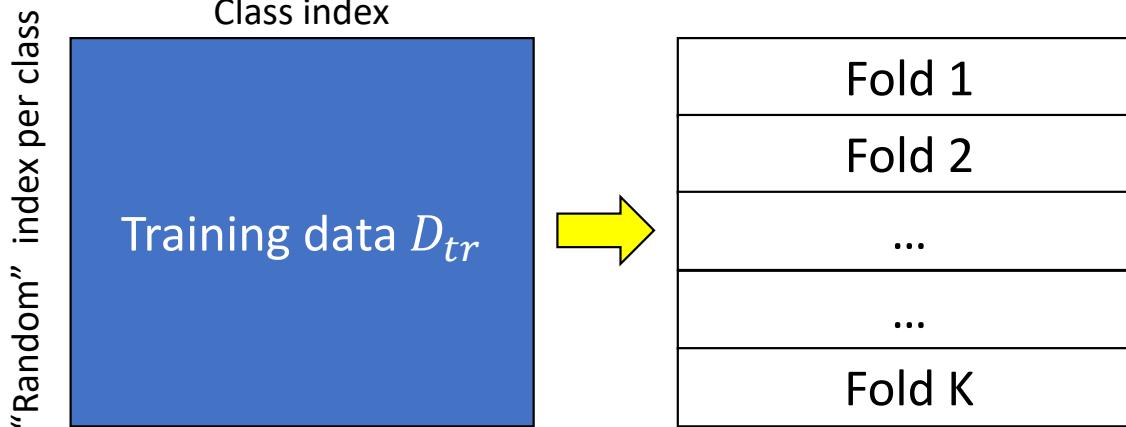


(Cross) validation

- Goal: choose the hyper-parameters
 - Ex: # of neighbors and distance in KNN, λ in ridge regression, η in iterative optimization

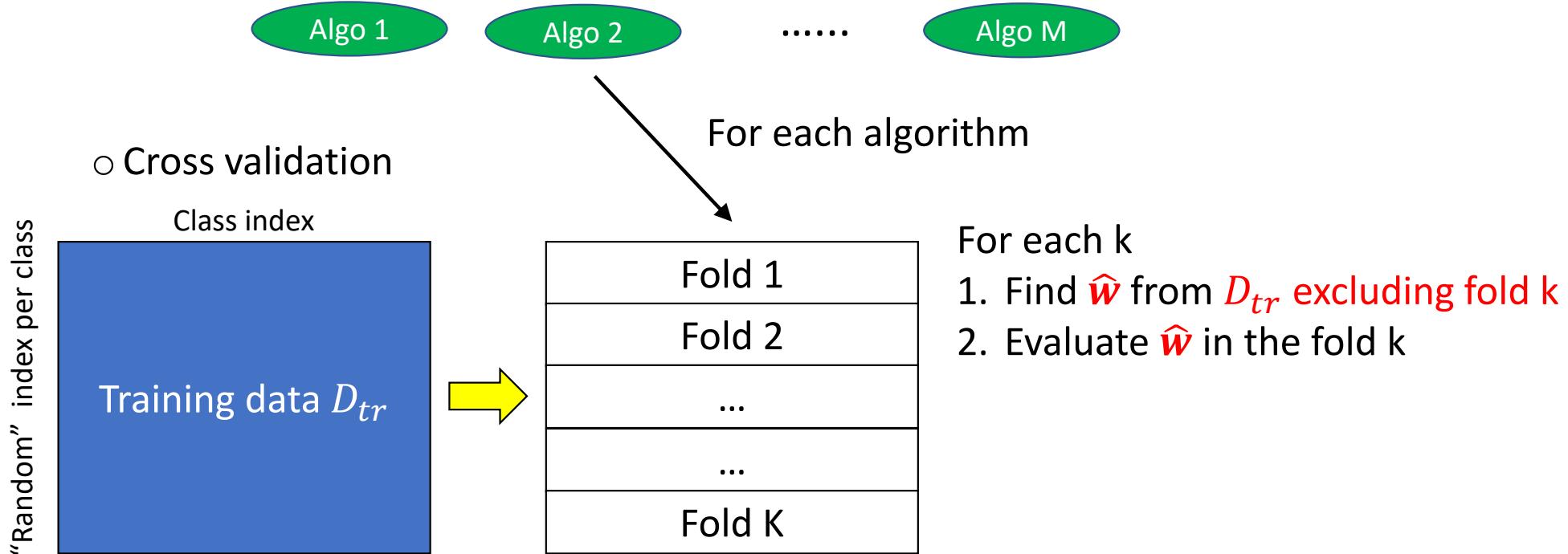


- Cross validation



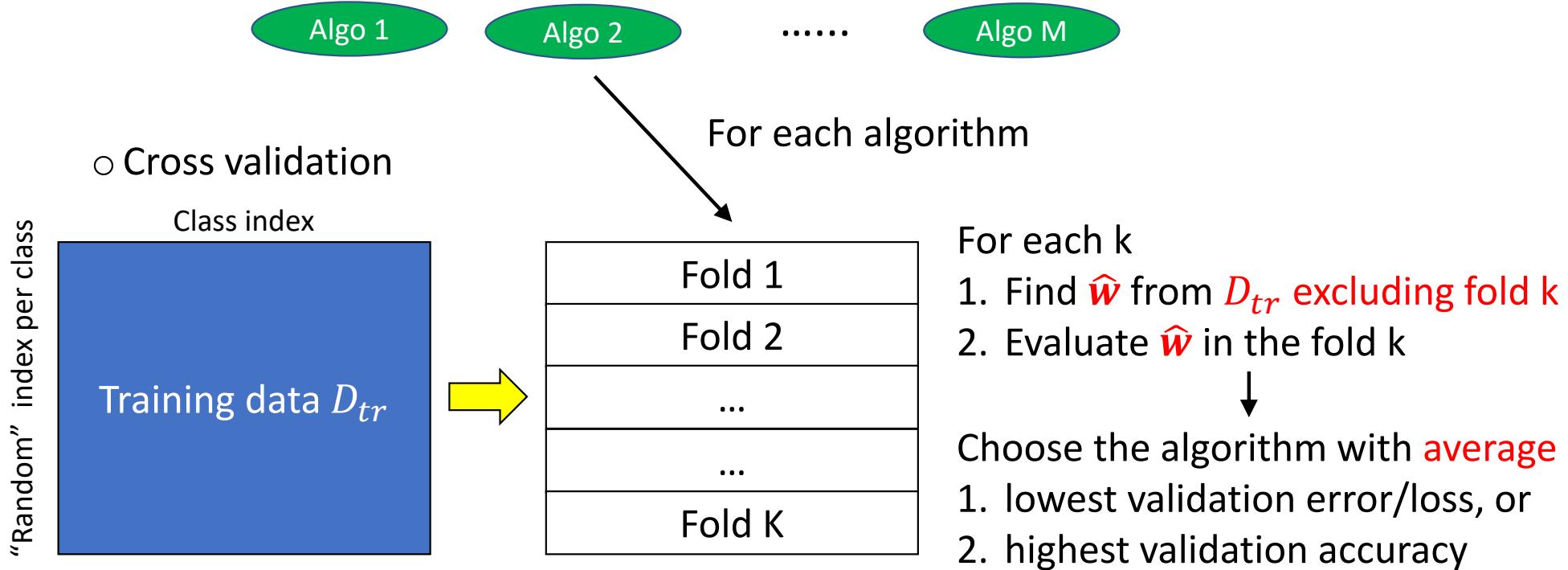
(Cross) validation

- Goal: choose the hyper-parameters
 - Ex:# of neighbors and distance in KNN, λ in ridge regression, η in iterative optimization



(Cross) validation

- Goal: choose the hyper-parameters
 - Ex:# of neighbors and distance in KNN, λ in ridge regression, η in iterative optimization



CSE 5523: Generative Learning



THE OHIO STATE UNIVERSITY

HW

- HW-2 was released this by morning and is due on 10/1.

Today

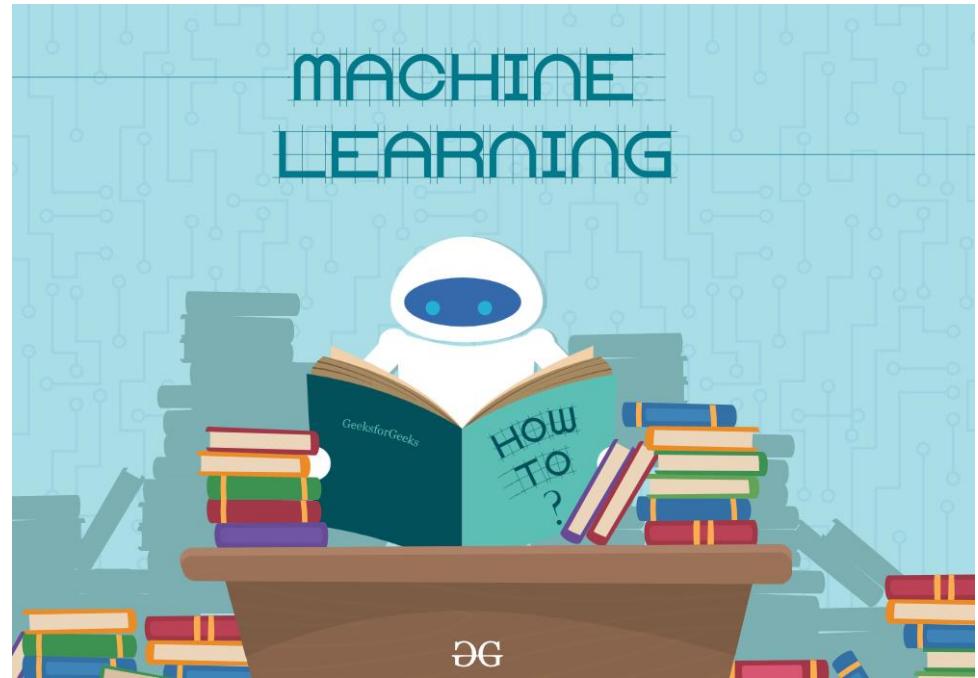
Review

Regularization, cross-validation

Gaussian discriminative analysis

Naïve Bayes

- Modeling
- Learning (informal)
- Learning (formal)



Estimating Probability from Data

- Data \mathcal{D} (observation/assignments): $D_{tr} = \{(x_n, y_n)\}_{n=1}^N$
 - Assumption: **IID** (or i.i.d.): **Independent** and **Identically Distributed**
- Goal: Fit a probability model p (interchangeable with P here) to data
 - $p(X, Y)$
 - $p(Y)p(X|Y)$: generative learning
 - $p(X)p(Y|X)$: discriminative learning
- Prediction
 - Given a test example x
 - $\hat{y} = \max_c p(Y = c|x)$; the Bayes Optimal Classifier

Estimating “Parameters of Probability” from Data

- **Maximum likelihood estimation (MLE)**

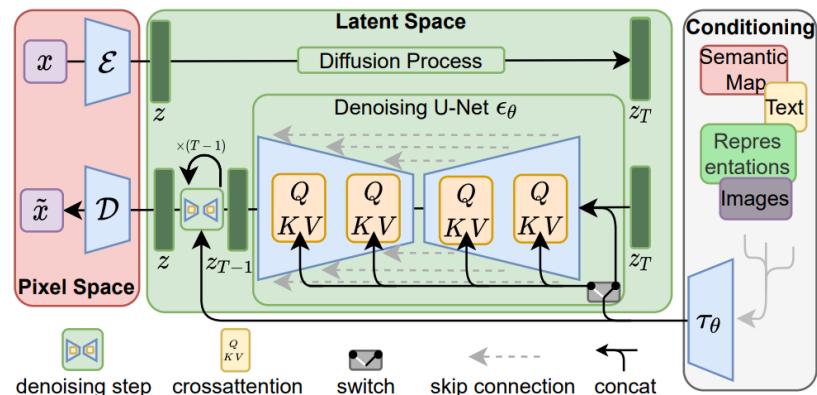
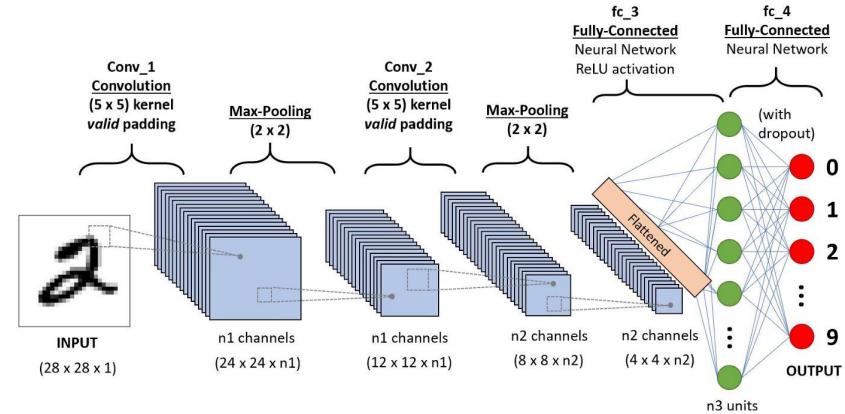
- $\widehat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\mathcal{D}; \boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N); \boldsymbol{\theta}) =$
 $\underset{\boldsymbol{\theta}}{\operatorname{argmax}} p((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N); \boldsymbol{\theta}) \times \dots p((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N); \boldsymbol{\theta}) =$
 $\underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_n p((\mathbf{x}_n, y_n); \boldsymbol{\theta})$

- **Maximum a posteriori estimation (MAP)**

- $\widehat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} P(\boldsymbol{\theta} | \mathcal{D}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\boldsymbol{\theta}) \times p((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) | \boldsymbol{\theta}) =$
 $\underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\boldsymbol{\theta}) \times \prod_n p((\mathbf{x}_n, y_n) | \boldsymbol{\theta})$

Don't get confused

- Discriminative vs. generative learning:
 - Difference in **modeling**
- MAP vs. MLE:
 - Difference in **the objective function**

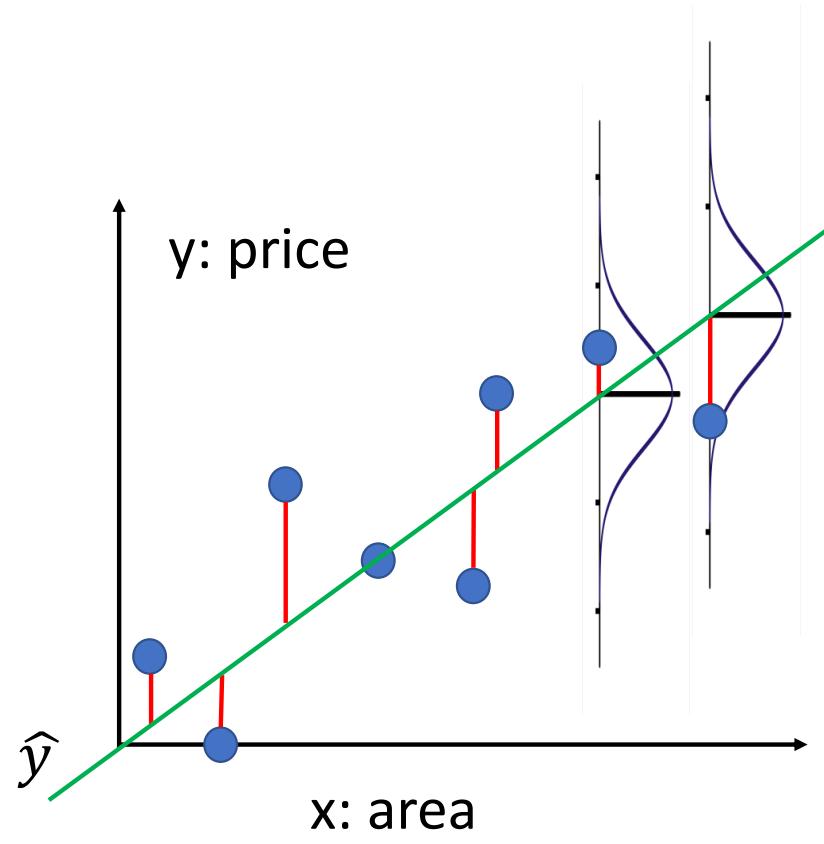


[Rombach et al., CVPR 2022]

Linear regression revisit (bias is absorbed)

- Modeling: $y|\mathbf{x} \sim \mathcal{N}(\mathbf{w}^T \mathbf{x}, \sigma^2)$
 - $p(y|\mathbf{x}; \mathbf{w}) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(y - \mathbf{w}^T \mathbf{x})^2 / 2\sigma^2}$

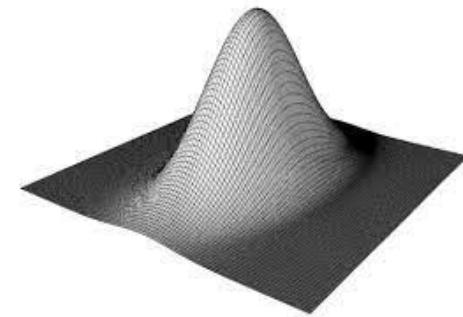
- MLE:
 - Given $D_{tr} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$
 - $\prod_{n=1}^N p(y_n|\mathbf{x}_n; \mathbf{w})$
 - $\log \prod_{n=1}^N p(y_n|\mathbf{x}_n; \mathbf{w})$
 - $\sum_n -(y_n - \mathbf{w}^T \mathbf{x}_n)^2$
 - $\hat{\mathbf{w}} = \operatorname{argmax} \sum_n -(y_n - \mathbf{w}^T \mathbf{x}_n)^2 =$
 $\operatorname{argmin}_{\theta} \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n)^2$



Ridge regression revisit

- MAP:

$$\circ p(\mathbf{w}) = \frac{1}{(2\pi)^{D/2} \det(\tau^2 \mathbf{I})^{1/2}} e^{-\frac{\frac{1}{2}\|\mathbf{w}\|_2^2}{\tau^2}}$$



$$\circ \log p(D_{tr} | \mathbf{w}) p(\mathbf{w}) \propto \sum_n - (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \log p(\mathbf{w})$$

$$\circ \operatorname{argmax}_{\mathbf{w}} \sum_n - (y_n - \mathbf{w}^T \mathbf{x}_n)^2 - \frac{1}{2} \frac{\|\mathbf{w}\|_2^2}{\tau^2} = \operatorname{argmin}_{\mathbf{w}} \left\{ \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n)^2 \right\} + \frac{1}{2} \frac{\|\mathbf{w}\|_2^2}{\tau^2}$$

- Ridge regression includes an L2 regularizer (not dependent on data)

Today

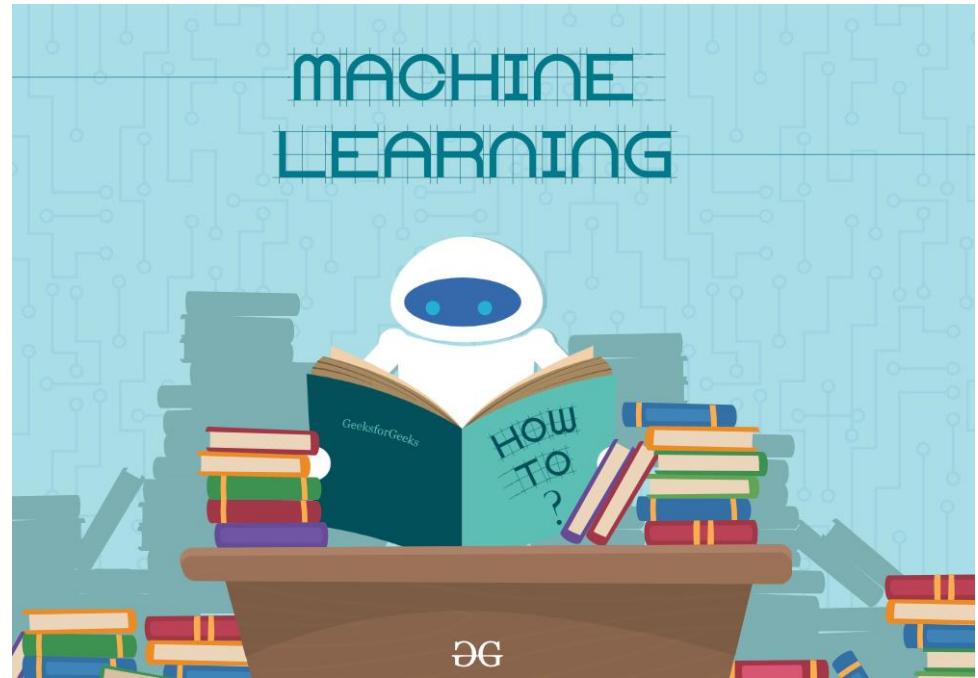
Review

Regularization, cross-validation

Gaussian discriminative analysis

Naïve Bayes

- Modeling
- Learning (informal)
- Learning (formal)



Regularization and hyper-parameters

- Ridge regression:

- $\underset{\mathbf{w}}{\operatorname{argmin}} \quad \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \lambda \|\mathbf{w}\|_2^2$

Regularization and hyper-parameters

- Ridge regression:

- $\underset{\mathbf{w}}{\operatorname{argmin}} \quad \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \lambda \|\mathbf{w}\|_2^2$

- The object function contains a training loss term and a regularization (e.g., belief) term

Regularization and hyper-parameters

- Ridge regression:

$$\underset{\mathbf{w}}{\operatorname{argmin}} \quad \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \lambda \|\mathbf{w}\|_2^2$$

- The **object function** contains a **training loss term** and a **regularization (e.g., belief) term**
- Model parameters (e.g., \mathbf{w}) is learned from minimizing the **training objective**

Regularization and hyper-parameters

- Ridge regression:

$$\underset{\mathbf{w}}{\operatorname{argmin}} \quad \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \lambda \|\mathbf{w}\|_2^2$$

- The **object function** contains a **training loss term** and a **regularization (e.g., belief) term**
- Model parameters (e.g., \mathbf{w}) is learned from minimizing the **training objective**
- How about hyper-parameters (e.g., λ) of the algorithm or hypothesis class?
 - Can you determine λ by minimizing the training objective?

Regularization and hyper-parameters

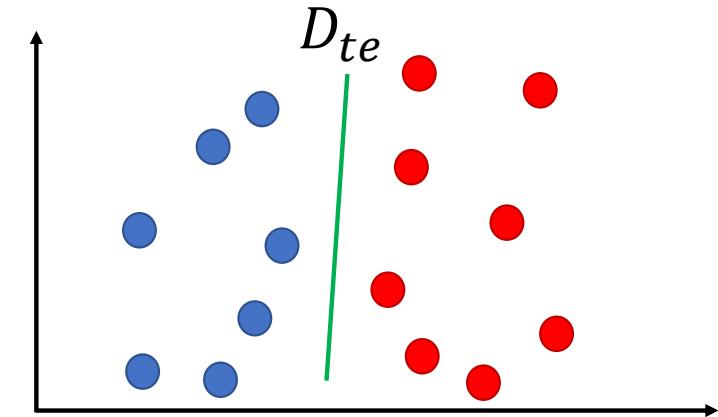
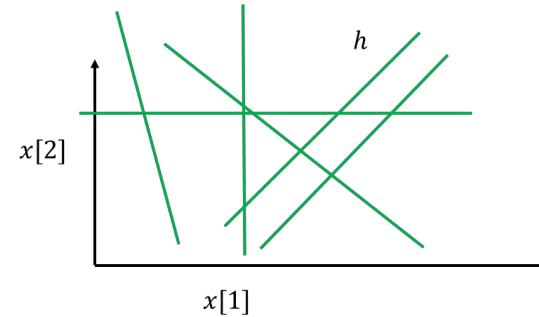
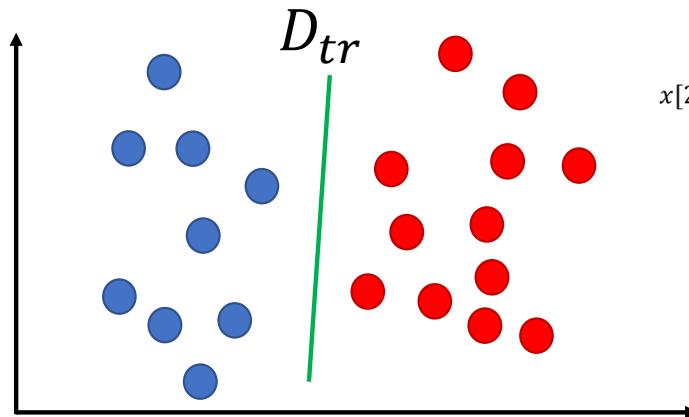
- Ridge regression:

$$\underset{\mathbf{w}}{\operatorname{argmin}} \quad \sum_n (y_n - \mathbf{w}^T \mathbf{x}_n)^2 + \lambda \|\mathbf{w}\|_2^2$$

- The **object function** contains a **training loss term** and a **regularization (e.g., belief) term**
- Model parameters (e.g., \mathbf{w}) is learned from minimizing the **training objective**
- How about hyper-parameters (e.g., λ) of the algorithm or hypothesis class?
 - Can you determine λ by minimizing the training objective?
- Caution: Don't be confused by notations; λ is not always used as parameters.

Machine learning

- A good learning algorithm + hypothesis class
 - Contain a hypothesis that can achieve low **training loss** and **test (validation) loss**
 - The algorithm can return such a hypothesis
 - We may view ridge regression as regularized loss (risk) minimization + linear regressors



- In practice, we choose the learning algorithm + hypothesis class by validation

(Cross) validation

- Goal: choose the hyper-parameters (ID of an algorithm + hypothesis class)
 - Ex:# of neighbors and distance in KNN, λ in ridge regression, η in iterative optimization

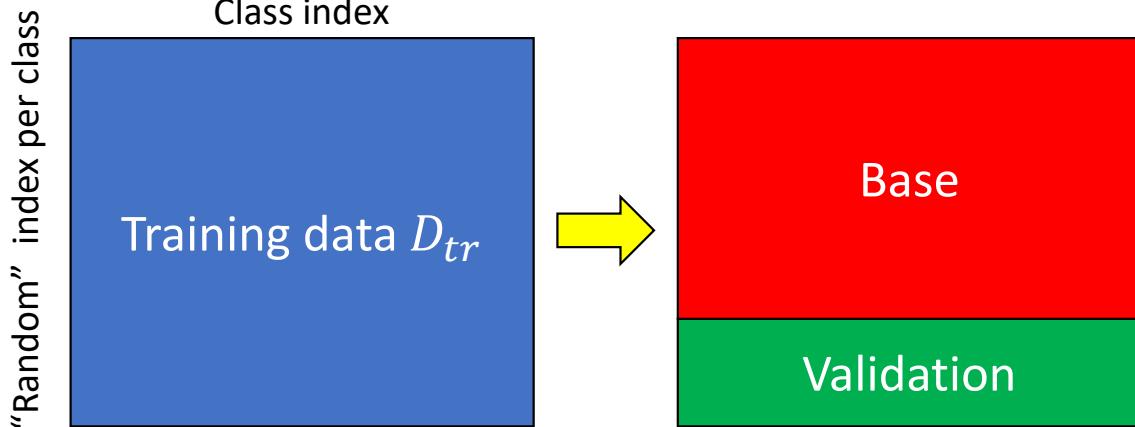


(Cross) validation

- Goal: choose the hyper-parameters
 - Ex:# of neighbors and distance in KNN, λ in ridge regression, η in iterative optimization

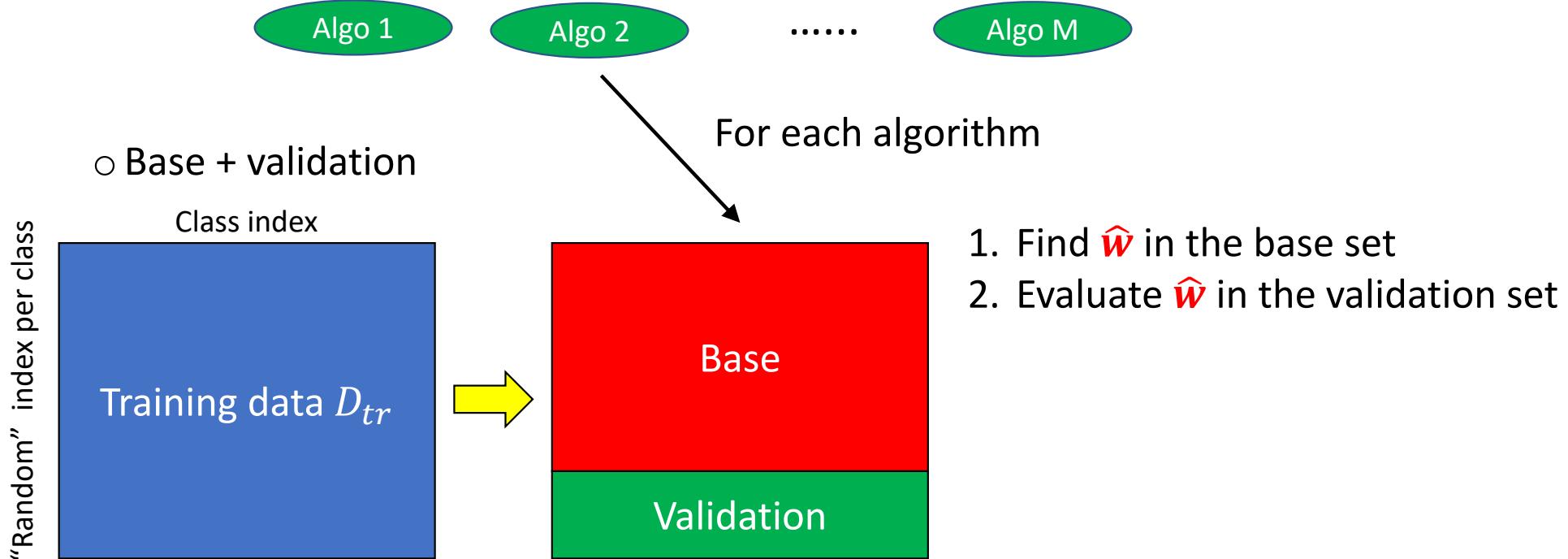


- Base + validation



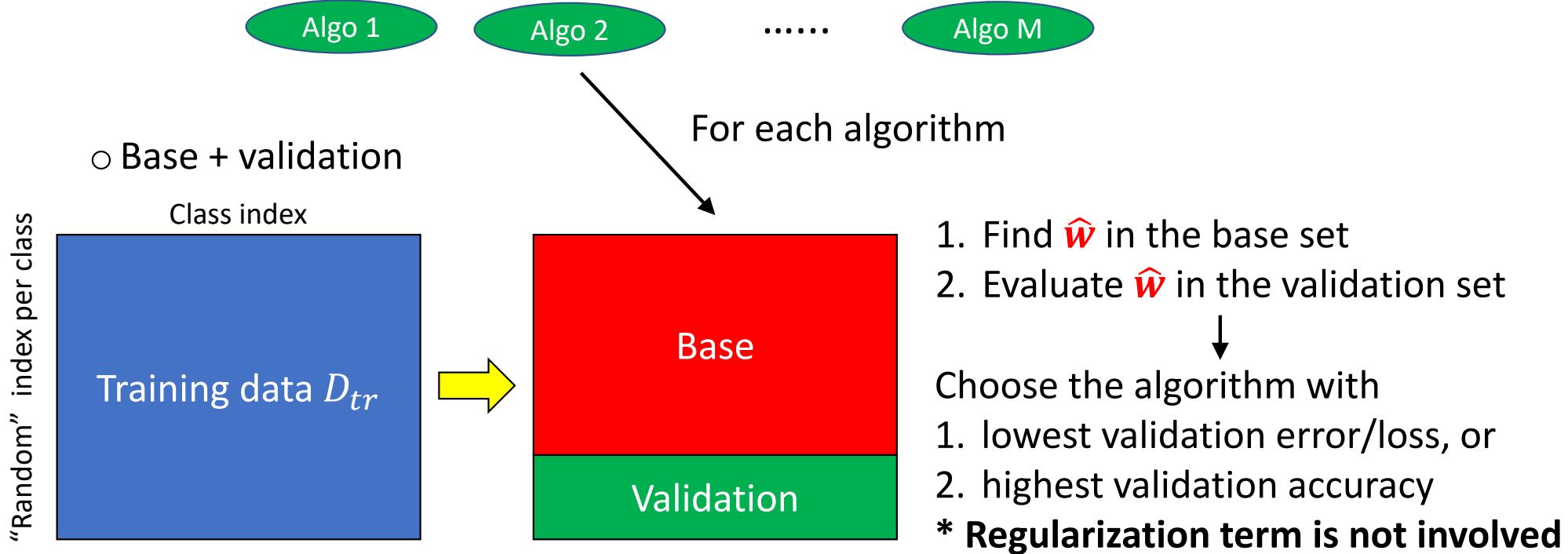
(Cross) validation

- Goal: choose the hyper-parameters
 - Ex:# of neighbors and distance in KNN, λ in ridge regression, η in iterative optimization



(Cross) validation

- Goal: choose the hyper-parameters
 - Ex:# of neighbors and distance in KNN, λ in ridge regression, η in iterative optimization

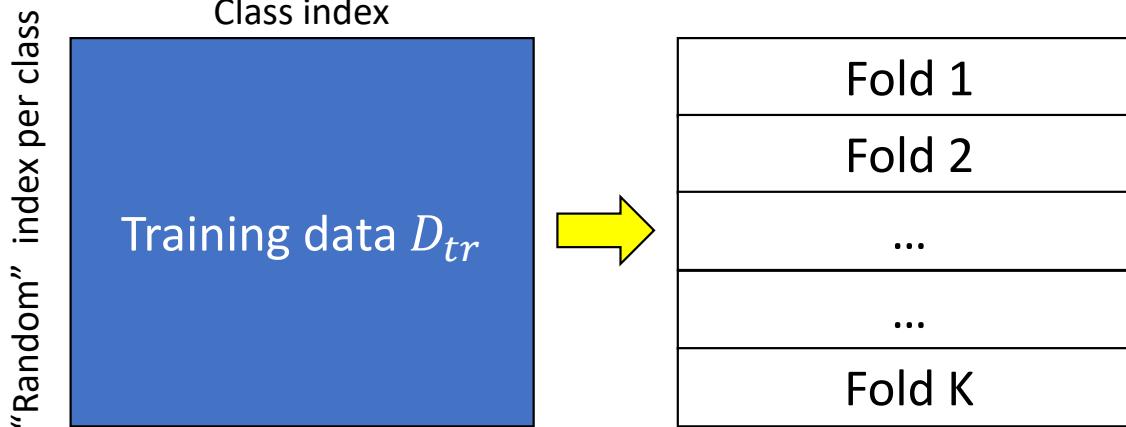


(Cross) validation

- Goal: choose the hyper-parameters
 - Ex: # of neighbors and distance in KNN, λ in ridge regression, η in iterative optimization

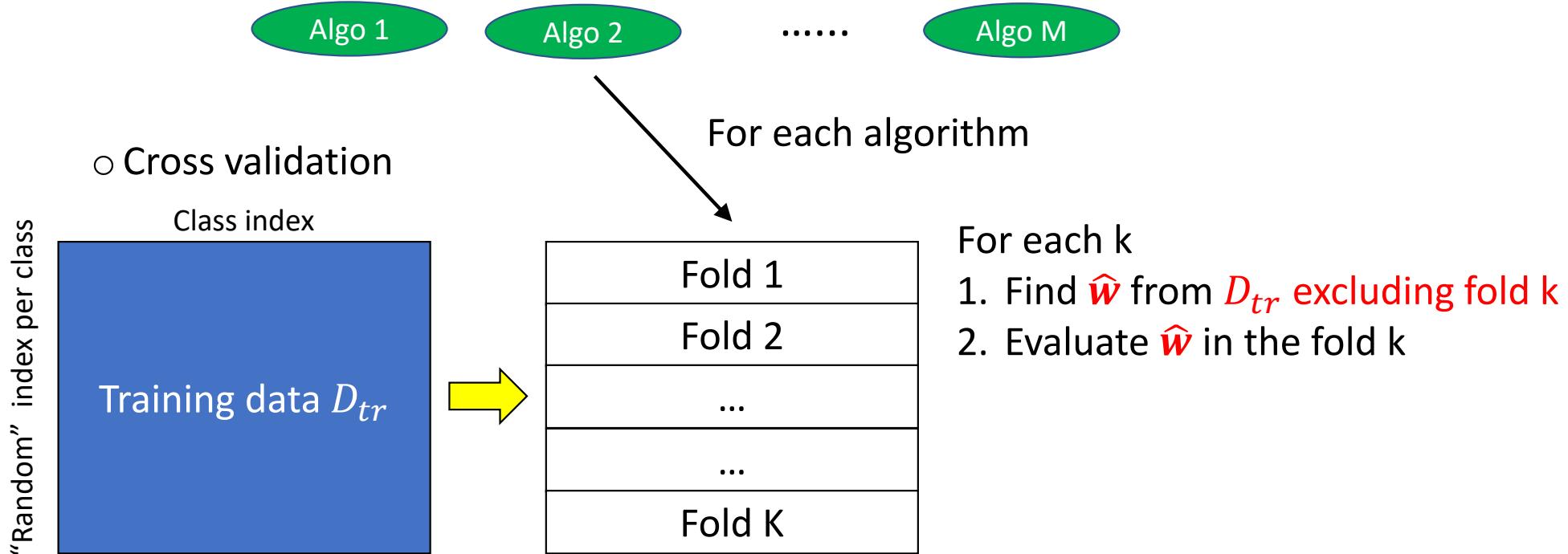


- Cross validation



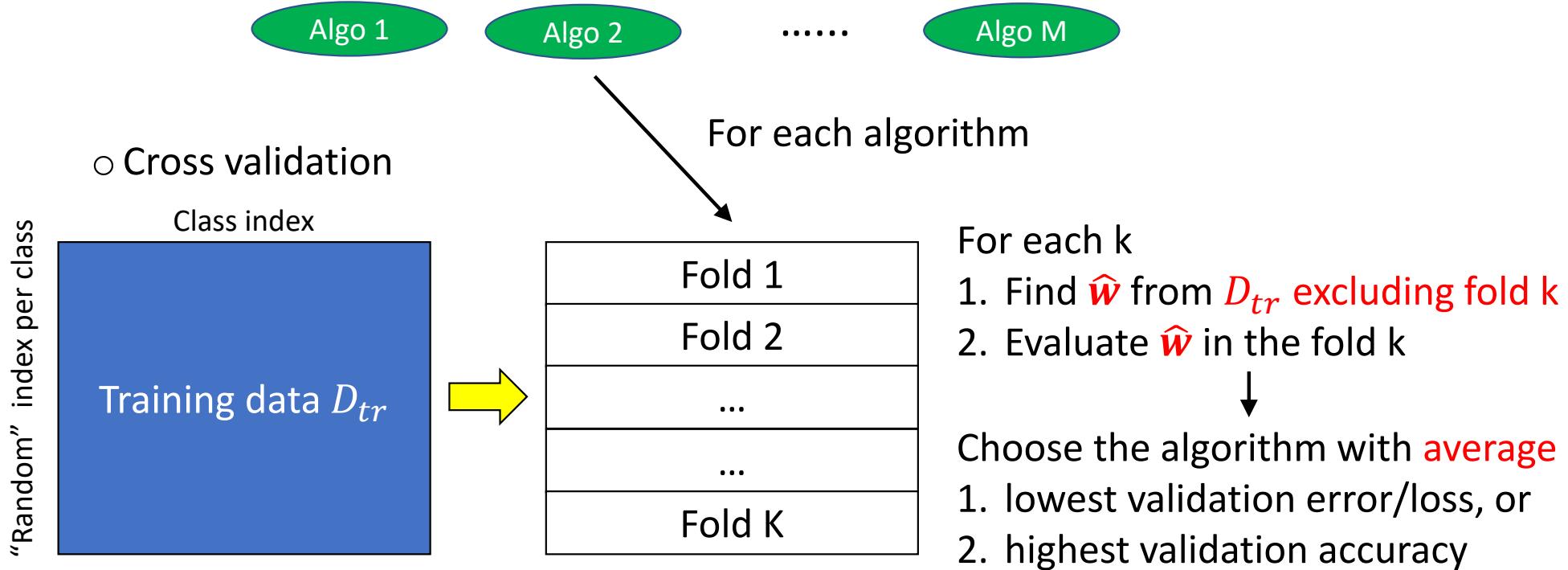
(Cross) validation

- Goal: choose the hyper-parameters
 - Ex:# of neighbors and distance in KNN, λ in ridge regression, η in iterative optimization



(Cross) validation

- Goal: choose the hyper-parameters
 - Ex:# of neighbors and distance in KNN, λ in ridge regression, η in iterative optimization



Today

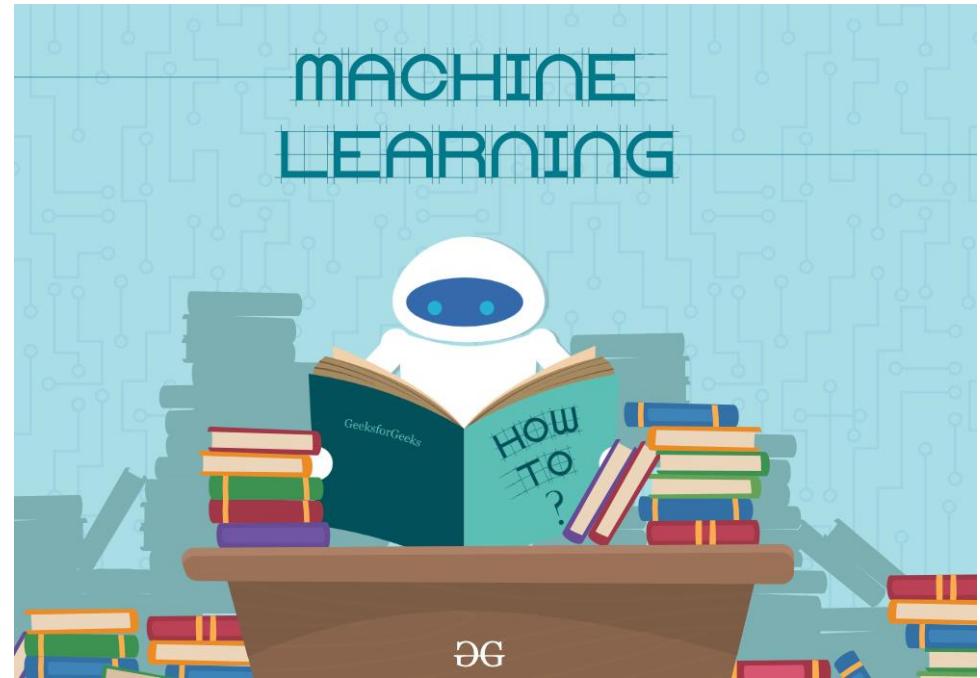
Review

Regularization, cross-validation

Gaussian discriminative analysis

Naïve Bayes

- Modeling
- Learning (informal)
- Learning (formal)



Gaussian discriminative analysis (GDA)

- Training data: $D_{tr} = \{(\mathbf{x}_n, y_n \in \{1, 2, \dots, C\})\}_{n=1}^N$
- To fit $p(X, Y) = p(Y)p(X|Y)$: so **generative** learning indeed
- Assumption:
 - $p(Y; \boldsymbol{\theta})$ is Bernoulli or categorical
 - $p(X|Y = y; \boldsymbol{\theta})$ is multivariate Gaussian
 - Sometimes, assuming shared covariance

Gaussian discriminative analysis (GDA)

- Training data: $D_{tr} = \{(\mathbf{x}_n, y_n \in \{1, 2, \dots, C\})\}_{n=1}^N$
- To fit $p(X, Y) = p(Y)p(X|Y)$: so **generative** learning indeed
- Assumption:
 - $p(Y; \boldsymbol{\theta})$ is Bernoulli or categorical
 - $p(X|Y = y; \boldsymbol{\theta})$ is multivariate Gaussian
 - Sometimes, assuming shared covariance
- Prediction: $\operatorname{argmax}_y \log p(y|\mathbf{x}; \hat{\boldsymbol{\theta}}) = \operatorname{argmax}_y \log p(y; \hat{\boldsymbol{\theta}}) + \log p(\mathbf{x}|y; \hat{\boldsymbol{\theta}})$

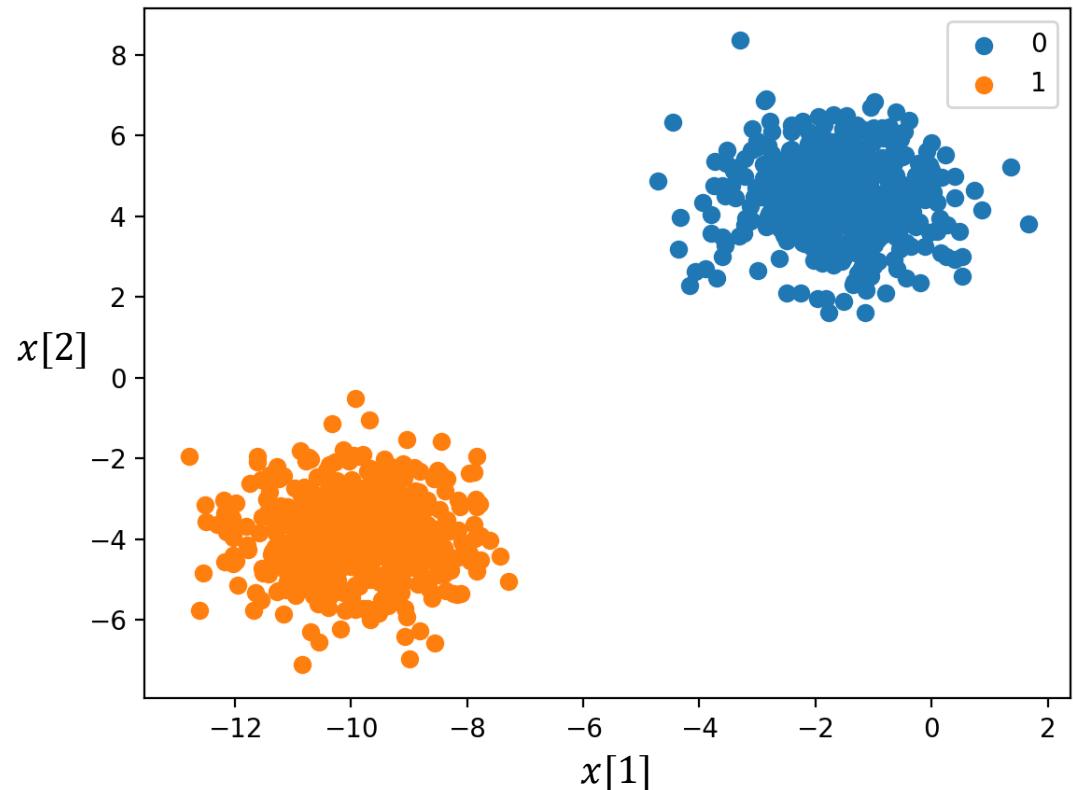
Given \mathbf{x} , which y is the most likely?

Which y is the most likely?

Which y most likely generates \mathbf{x} ?

Gaussian discriminative analysis

$$p(X, Y) = p(Y)p(X|Y)$$

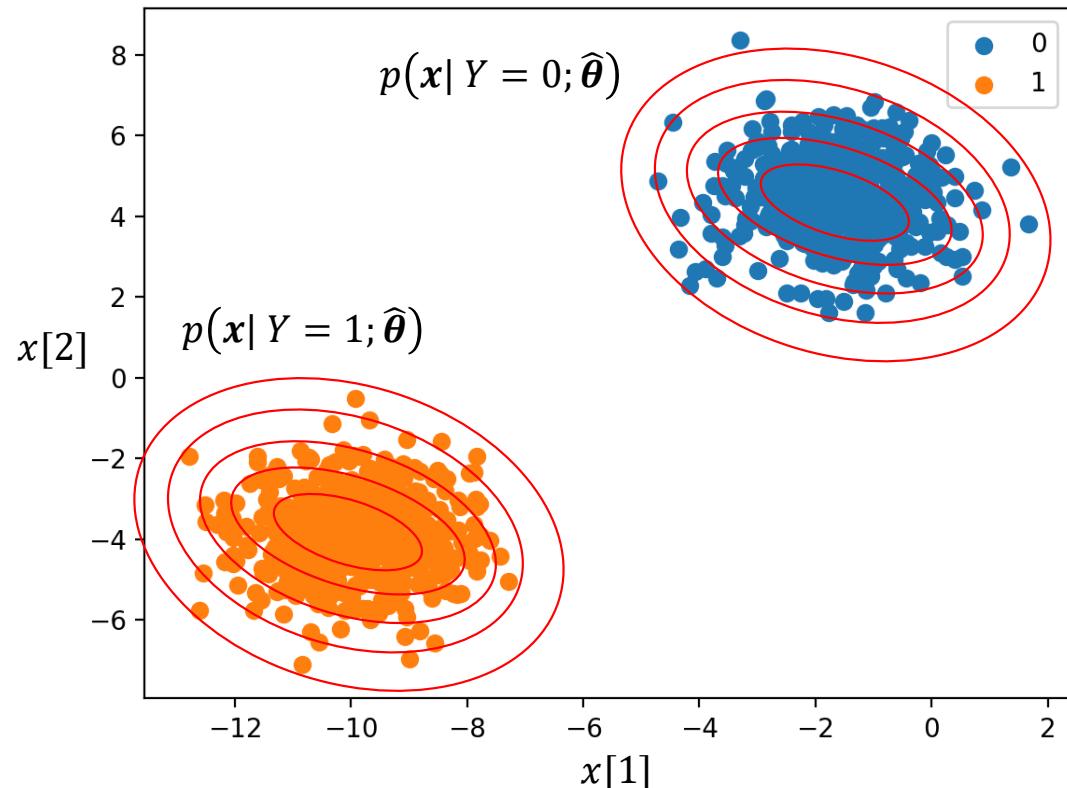


Gaussian discriminative analysis

- Intuition:

- $p(Y = 1; \hat{\theta}) =$
- $p(Y = 0; \hat{\theta}) =$
- $p(x| Y = 1; \hat{\theta}) =$
- $p(x| Y = 0; \hat{\theta}) =$
- θ or $\hat{\theta}$ is the collection of all model parameters

$$p(X, Y) = p(Y)p(X|Y)$$

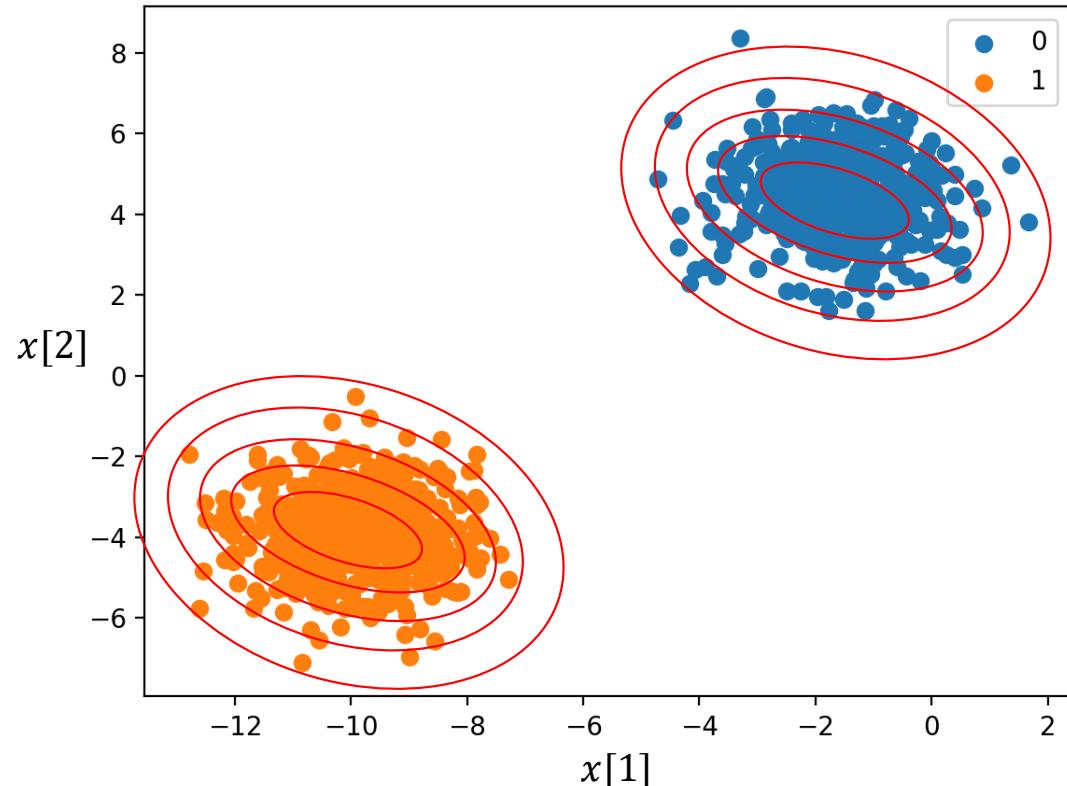


Gaussian discriminative analysis

- Intuition:

- $p(Y = 1; \theta) = \lambda = \frac{N_1}{N_1 + N_0}$
- $p(Y = 0; \theta) = \frac{N_0}{N_1 + N_0}$
- $p(x|Y = 1; \theta)$ by μ_1 and Σ_1 from all orange points
- $p(x|Y = 0; \theta)$ by μ_0 and Σ_0 from all blue points

$$p(X, Y) = p(Y)p(X|Y)$$



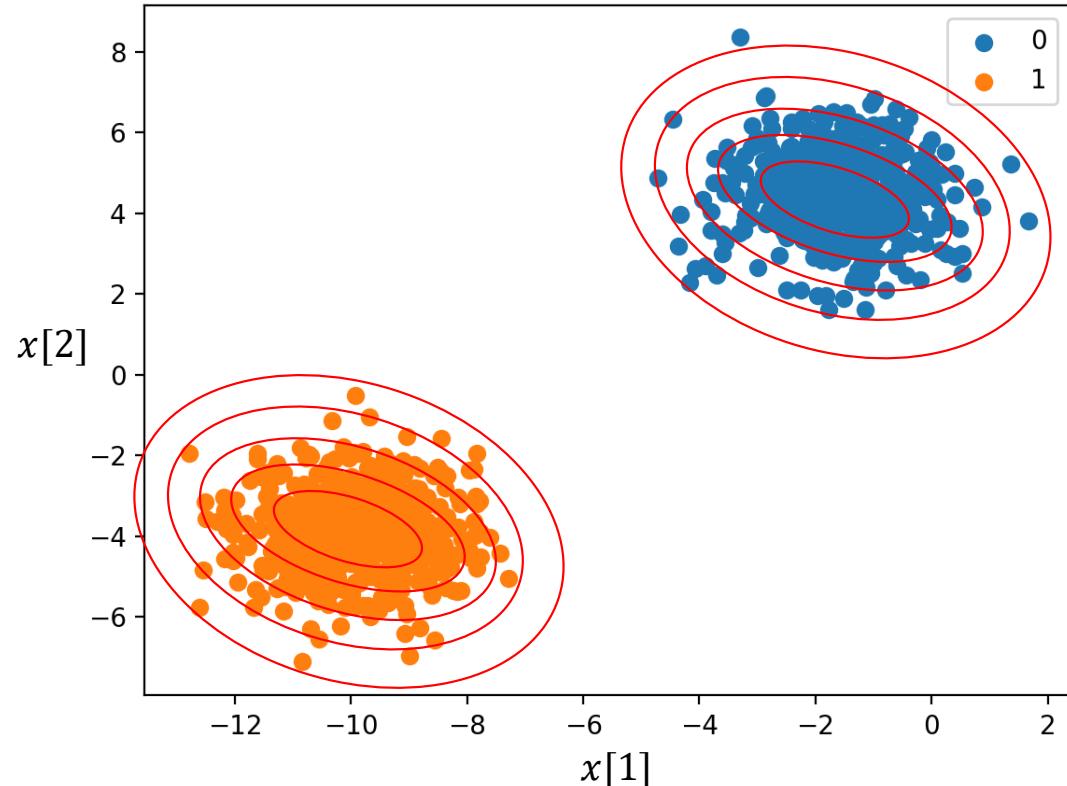
Gaussian discriminative analysis

- Intuition:

- $p(Y = 1; \theta) = \lambda = \frac{N_1}{N_1 + N_0}$
- $p(Y = 0; \theta) = \frac{N_0}{N_1 + N_0}$
- $p(x|Y = 1; \theta)$ by μ_1 and Σ_1 from all orange points
- $p(x|Y = 0; \theta)$ by μ_0 and Σ_0 from all blue points

- Shared Σ (i.e., $\Sigma_1 = \Sigma_0$):

$$\Sigma = \frac{N_1 \Sigma_1 + N_0 \Sigma_0}{N_1 + N_0}$$



Gaussian discriminative analysis

- Shared Σ (i.e., $\Sigma_1 = \Sigma_0$):

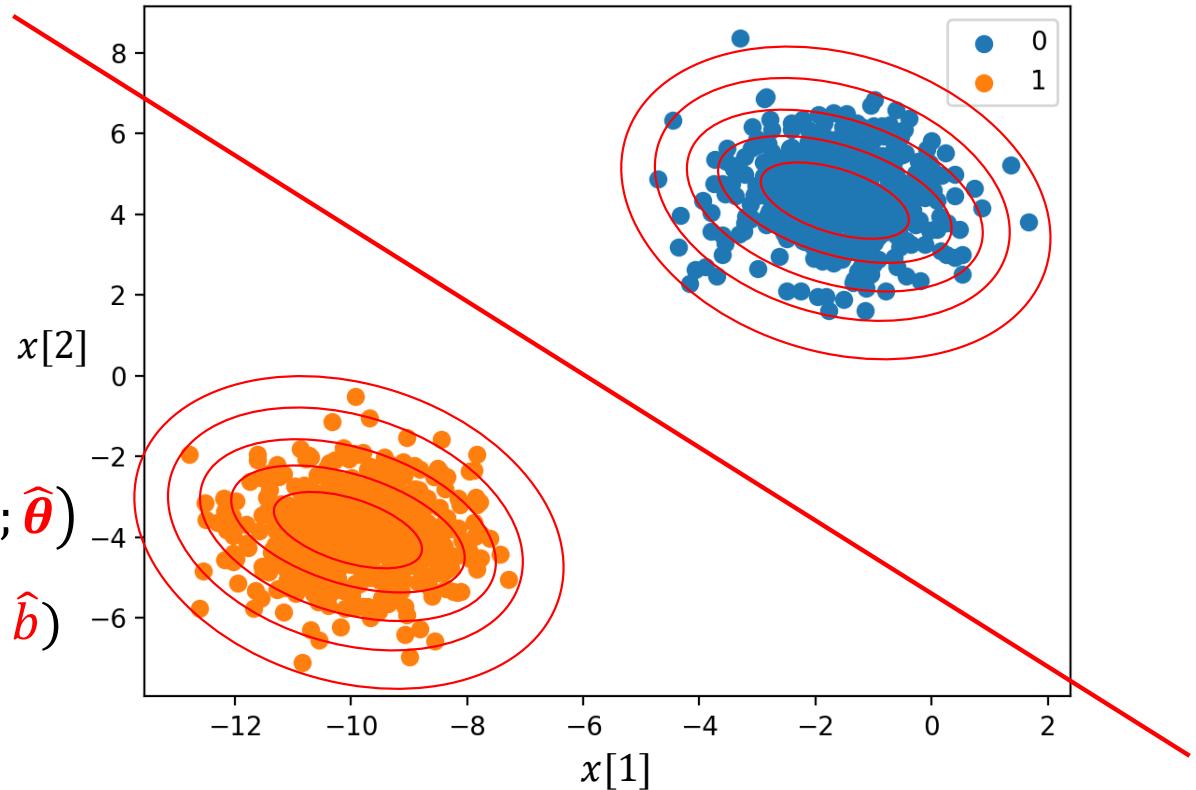
$$\Sigma = \frac{N_1 \Sigma_1 + N_0 \Sigma_0}{N_1 + N_0}$$

- Prediction

- $\underset{y}{\operatorname{argmax}} \log p(y; \hat{\theta}) + \log p(x|y; \hat{\theta})$

- Can be rewritten as $\operatorname{sign}(\hat{w}^T x + \hat{b})$

- Output 1 means class 1
 - Output -1 means class 0



Gaussian discriminative analysis

- MLE on the joint probability

$$\begin{aligned} \log \prod_{i=1}^N p(\mathbf{x}_i, y_i; \boldsymbol{\theta}) &= \sum_{i=1}^N \log p(\mathbf{x}_i, y_i; \boldsymbol{\theta}) = \sum_{i=1}^N \log p(\mathbf{x}_i|y_i; \boldsymbol{\theta})p(y_i; \boldsymbol{\theta}) \\ &= \sum_{i=1}^N \log p(\mathbf{x}_i|y_i; \boldsymbol{\theta}) + \sum_{i=1}^N \log p(y_i; \boldsymbol{\theta}) \end{aligned}$$

Gaussian for $y_i = 1$ or $y_i = 0$

Bernoulli, do not share parameters with Gaussian

- Each term is affected by some “separate” parameters
 - We can then estimate each parameter based on a certain set of data

Gaussian discriminative analysis

- MLE on the joint probability

$$\begin{aligned} \log \prod_{i=1}^N p(\mathbf{x}_i, y_i; \boldsymbol{\theta}) &= \sum_{i=1}^N \log p(\mathbf{x}_i, y_i; \boldsymbol{\theta}) = \sum_{i=1}^N \log p(\mathbf{x}_i|y_i; \boldsymbol{\theta})p(y_i; \boldsymbol{\theta}) \\ &= \sum_{i=1}^N \log p(\mathbf{x}_i|y_i; \boldsymbol{\theta}) + \sum_{i=1}^N \log p(y_i; \boldsymbol{\theta}) \end{aligned}$$

Gaussian for $y_i = 1$ or $y_i = 0$

Bernoulli, do not share parameters with Gaussian

- $\underset{\lambda}{\operatorname{argmax}} \sum_{i=1}^N \log p(\mathbf{x}_i|y_i) + \sum_{i=1}^N \log p(y_i) = \underset{\lambda}{\operatorname{argmax}} \sum_{i=1}^N \log p(y_i)$
 - Where did you see this?

Gaussian discriminative analysis

- MLE on the joint probability

$$\begin{aligned} \log \prod_{i=1}^N p(\mathbf{x}_i, y_i; \boldsymbol{\theta}) &= \sum_{i=1}^N \log p(\mathbf{x}_i, y_i; \boldsymbol{\theta}) = \sum_{i=1}^N \log p(\mathbf{x}_i|y_i; \boldsymbol{\theta})p(y_i; \boldsymbol{\theta}) \\ &= \sum_{i=1}^{\textcolor{red}{N}} \log p(\mathbf{x}_i|y_i; \boldsymbol{\theta}) + \sum_{i=1}^{\textcolor{blue}{N}} \log p(y_i; \boldsymbol{\theta}) \end{aligned}$$

Gaussian for $y_i = 1$ or $y_i = 0$

Bernoulli, do not share parameters with Gaussian

- $\underset{\Sigma_1}{\operatorname{argmax}} \sum_{i=1}^N \log p(\mathbf{x}_i|y_i) + \sum_{i=1}^N \log p(y_i) = \underset{\Sigma_1}{\operatorname{argmax}} \sum_{i \in \{y_i=1\}} \log p(\mathbf{x}_i|y_i)$

Is this model powerful?

- Not really if the data (feature vectors) per class is not Gaussian
- Usually powerful when paired with deep learning features
- Efficient model updates:
 - Let $\boldsymbol{\theta} = \{\lambda, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2\}$
 - Let $\widehat{\boldsymbol{\theta}}^{(N)}$ be the parameters estimated from $D_{tr} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$
 - If now there is a new training example $(\mathbf{x}_{N+1}, y_{N+1})$ coming in
 - To estimate $\widehat{\boldsymbol{\theta}}^{(N+1)}$ from $D_{tr} + \{(\mathbf{x}_{N+1}, y_{N+1})\}$, we don't need to keep D_{tr}
 - Instead, we only need to keep $\widehat{\boldsymbol{\theta}}^{(N)}$ and some other information like N_1, N_0
 - $\boldsymbol{\theta} + \{N_1, N_0\}$ can be viewed as the sufficient statistics of D_{tr} for GDA
 - That is, knowing $\boldsymbol{\theta} + \{N_1, N_0\}$ of D_{tr} is sufficient for building up the GDA models

Today

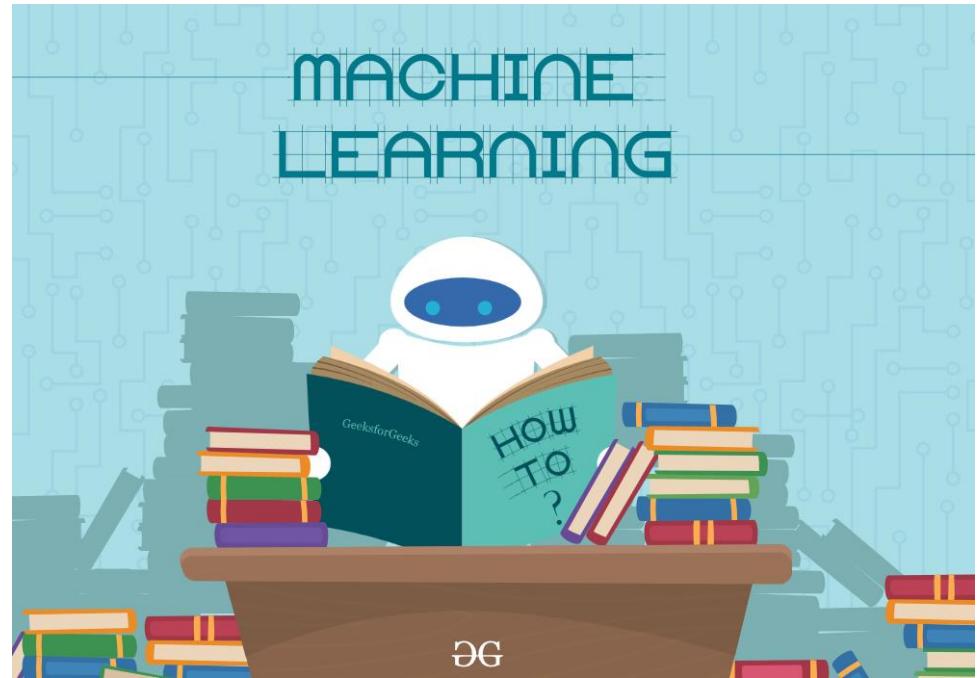
Review

Regularization, cross-validation

Gaussian discriminative analysis

Naïve Bayes

- Modeling
- Learning (informal)
- Learning (formal)



Naïve Bayes

- **Training data:** $D_{tr} = \{ (x_i \in \mathbb{R}^D, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$
 - Each dimension of x_i can be binary, categorical, count, or numerical (real)
- **Goal:** construct $p(Y = c|x)$ for $\hat{y} = \max_c p(Y = c|x)$

Naïve Bayes

- **Training data:** $D_{tr} = \{ (x_i \in \mathbb{R}^D, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$
 - Each dimension of x_i can be binary, categorical, count, or numerical (real)
- **Goal:** construct $p(Y = c|x)$ for $\hat{y} = \max_c p(Y = c|x)$
- **How? Bayes' rules**
 - $p(Y = c|x) = \frac{P(x|Y=c)p(Y=c)}{P(x)} \propto p(x|Y=c)p(Y=c)$
 - $\operatorname{argmax}_c p(Y = c|x) = \operatorname{argmax}_c p(x|Y=c)p(Y=c)$
- **Key:** what probability models for $p(x|Y = c; \theta)$ and $p(Y = c; \theta)$?

Naïve Bayes: modeling

- $p(Y = c; \theta)$: Bernoulli or categorical
- $p(x|Y = c; \theta)$: **conditional independence assumption:**
 - $p(\textcolor{red}{x}|Y = c) = p(x[1], x[2], x[3], \dots, x[D] | Y = c)$
 - $= p(x[1] | Y = c)p(x[2] | Y = c) \dots \dots p(x[D] | Y = c) = \prod_{d=1}^D p(x[d] | Y = c)$

Naïve Bayes: modeling

- $p(Y = c; \theta)$: Bernoulli or categorical
- $p(x|Y = c; \theta)$: **conditional independence assumption:**
 - $p(x|Y = c) = p(x[1], x[2], x[3], \dots, x[D] | Y = c)$
 $= p(x[1] | Y = c)p(x[2] | Y = c) \dots \dots p(x[D] | Y = c) = \prod_{d=1}^D p(x[d] | Y = c)$
 - Given the class label, the probability of observing the value of one dimension is conditionally independent from the values of the other dimensions
 - Ex: Spam filtering
 - x : an email; $x[d]$: the d -th word or the d -th unique word; $c \in \{\text{spam}, \text{not spam}\}$
 - Clearly, this assumption is usually not true, but the model works well in practice.
 - "All models are wrong, but some are useful."

Naïve Bayes: data generation and inference

- Example:

- Every time, sample a coin (from two “types”)
- Flip the sampled coin and measure the weight



| Instance ID: i | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|-------------------|-----|-----|-----|-----|-----|-----|-----|-----|
| Coin | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| Head (1)/Tail (0) | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 0 |
| Weight | 2.2 | 2.3 | 2.5 | 1.7 | 1.9 | 2.1 | 1.5 | 2.7 |

y_i
 $x_i[1]$
 $x_i[2]$

- Given the type of coin, observing head/tail is conditionally independent from the weight

- **Classification problem:** Now given a new observation x (e.g., “head” and “weight = 2.4”), what type of coin (0 or 1) is it from?

Naïve Bayes: data generation and inference

- Example:

- Every time, sample a coin (from two “types”)
- Flip the sampled coin and measure the weight



| Instance ID: i | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|-------------------|-----|-----|-----|-----|-----|-----|-----|-----|
| Coin | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| Head (1)/Tail (0) | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 0 |
| Weight | 2.2 | 2.3 | 2.5 | 1.7 | 1.9 | 2.1 | 1.5 | 2.7 |

y_i
 $x_i[1]$
 $x_i[2]$

- Given the type of coin, observing head/tail is conditionally independent from the weight

- $p(Y = 1) \times p(x[1] = \text{head}|Y = 1) \times p(x[2] = 2.4|Y = 1)$ vs.
- $p(Y = 0) \times p(x[1] = \text{head}|Y = 0) \times p(x[2] = 2.4|Y = 0)$

Today

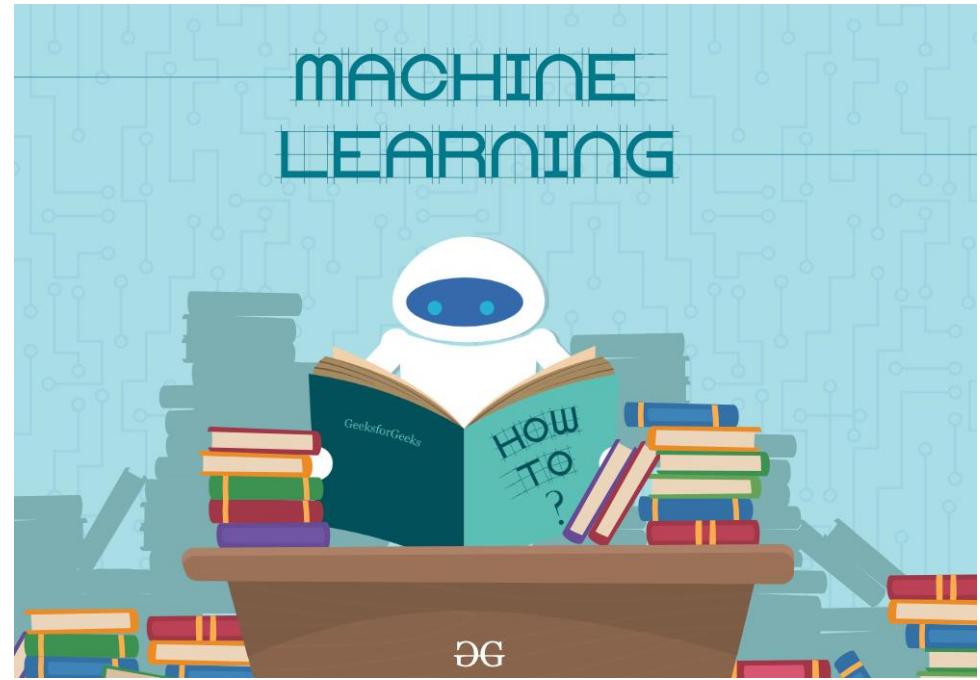
Review

Regularization, cross-validation

Gaussian discriminative analysis

Naïve Bayes

- Modeling
- Learning (informal)
- Learning (formal)



Naïve Bayes: another example

x is 3-dimensional car profile; $D_{tr} = \{ (x_i, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$

Maintained?

Car brand?

Mileage?

$$\begin{aligned}(x|Y = c) \\ &= p(x[1] | Y = c) \\ &\times p(x[2] | Y = c) \\ &\times p(x[3] | Y = c)\end{aligned}$$

Naïve Bayes: another example

x is 3-dimensional car profile; $D_{tr} = \{ (x_i, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$

| |
|-------------|
| Maintained? |
| Car brand? |
| Mileage? |



Bernoulli distribution: $p(x[d] = 1 | Y = c) = \lambda_{d,c}$

$$\lambda_{d,c} = \frac{\sum_{i=1}^N \mathbf{1}[y_i = c] \mathbf{1}[x_i[d] = 1]}{\sum_{i=1}^N \mathbf{1}[y_i = c]}$$

Each pair of (dimension, class) has a specific coin

The generative process:

1. Sample one class c
2. Flip the corresponding coin with $\lambda_{d,c}$

Naïve Bayes: another example

x is 3-dimensional car profile; $D_{tr} = \{ (x_i, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$

| |
|-------------|
| Maintained? |
| Car brand? |
| Mileage? |



Categorical distribution: $p(x[d] = m | Y = c) = \lambda_{d,c}[m]$

$$\lambda_{d,c}[m] = \frac{\sum_{i=1}^N \mathbf{1}[y_i = c] \mathbf{1}[x_i[d] = m]}{\sum_{i=1}^N \mathbf{1}[y_i = c]}$$

category m

Each pair of (dimension, class) has a specific dice

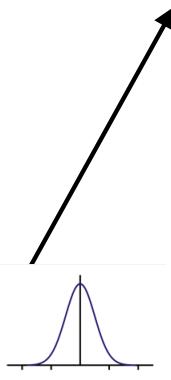
The generative process:

1. Sample one class c
2. Throw the corresponding dice with $\lambda_{d,c}$

Naïve Bayes: another example

x is 3-dimensional car profile; $D_{tr} = \{ (x_i, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$

| |
|-------------|
| Maintained? |
| Car brand? |
| Mileage? |



Gaussian distribution: $p(x[d]|Y = c) = \mathcal{N}(\mu_{d,c}, \sigma_{d,c}^2)$

$$\mu_{d,c} = \frac{1}{N_c} \sum_{i=1}^N \mathbf{1}[y_i = c] x[d]$$

$$\sigma_{d,c}^2 = \frac{1}{N_c} \sum_{i=1}^N \mathbf{1}[y_i = c] (x[d] - \mu_{d,c})^2$$

Each pair of (dimension, class) has a specific Gaussian
The generative process:

1. Sample one class c
2. Sample from the corresponding Gaussian with $\mu_{d,c}$ and $\sigma_{d,c}^2$

Today

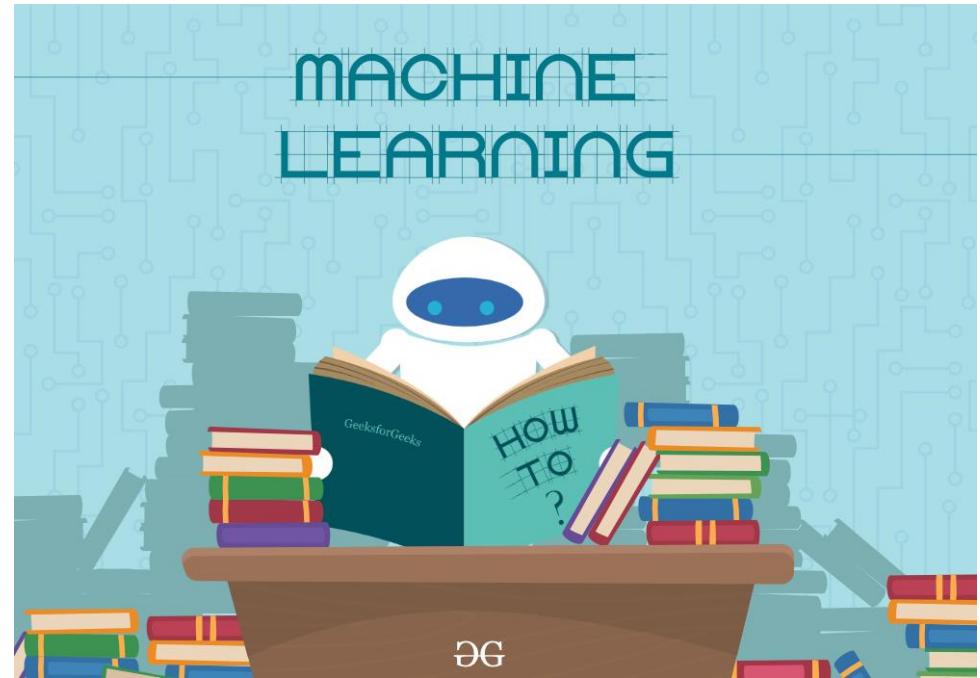
Review

Regularization, cross-validation

Gaussian discriminative analysis

Naïve Bayes

- Modeling
- Learning (informal)
- Learning (formal)



How to learn $p(\mathbf{x}[d] \mid Y = c)$ and $p(Y = c)$?

- Given training data
 - $\mathcal{D} = \{(\mathbf{x}_i, y_i \in \{1, 2, \dots, C\})\}_{i=1}^N$
- Applying the MLE criterion
 - $\underset{\text{parameters}}{\operatorname{argmax}} p(\mathcal{D}; \text{parameters})$
- $p(\mathcal{D}; \text{parameters}) = p(\{\mathbf{x}_1, y_1\}, \dots, \{\mathbf{x}_N, y_N\}) = \prod_{i=1}^N p(\mathbf{x}_i, y_i)$
 - IID observations

Naïve Bayes: learning

- MLE:
 - Maximize the “joint probability” of the **IID** data instances

$$\log p(\{x_1, y_1\}, \dots) = \log \prod_{i=1}^N p(x_i, y_i) = \sum_{i=1}^N \log p(x_i, y_i) = \sum_{i=1}^N \log p(x_i|y_i)p(y_i)$$

Naïve Bayes: learning

- MLE:

- Maximize the “joint probability” of the **IID** data instances
- Separately optimize each probability model: each has independent parameters and can have different models!

$$\sum_{d=1}^D \sum_{i=1}^N \log p(\mathbf{x}_i[d] | y_i) + \sum_{i=1}^N \log p(y_i)$$
$$= \boxed{\sum_{i=1}^N \log p(\mathbf{x}_i[1] | y_i)} + \dots \dots + \boxed{\sum_{i=1}^N \log p(\mathbf{x}_i[D] | y_i)} + \boxed{\sum_{i=1}^N \log p(y_i)}$$

Categorical, Bernoulli, or 1-D Gaussian for each class c

Categorical or Bernoulli

Naïve Bayes: learning

- MLE: for each conditional term $\sum_{i=1}^N \log p(\mathbf{x}_i[d] | y_i)$, we must build a conditional probability model for each class c
 - $p(\mathbf{x}[d] | Y = 1), p(\mathbf{x}[d] | Y = 2), \dots \dots, p(\mathbf{x}[d] | Y = C)$
- How?
 - Look at data instances of class c
 - Pick one feature dimension d
 - Perform MLE over $\mathbf{x}_i[d]$ whose $y_i = c$ (i.e., $\{(\mathbf{x}_i, y_i = c)\}$) to build $p(x[d] | Y = c)$
 - If $x[d]$ is continuous (e.g., weight), use a Gaussian
 - If $x[d]$ is binary, use Bernoulli; if $x[d]$ is categorical, use categorical
 - Of course, other distributions are applicable!

Naïve Bayes: learning

- Example:

| Instance ID: i | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|-------------------|-----|-----|-----|-----|-----|-----|-----|-----|
| Coin | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| Head (1)/Tail (0) | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 0 |
| Weight | 2.2 | 2.3 | 2.5 | 1.7 | 1.9 | 2.1 | 1.5 | 2.7 |

y_i
 $x_i[1]$
 $x_i[2]$

- What to learn?

Naïve Bayes: learning

- Example:

| Instance ID: i | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|-------------------|-----|-----|-----|-----|-----|-----|-----|-----|
| Coin | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| Head (1)/Tail (0) | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 0 |
| Weight | 2.2 | 2.3 | 2.5 | 1.7 | 1.9 | 2.1 | 1.5 | 2.7 |

y_i
 $x_i[1]$
 $x_i[2]$

- What to learn?
 - $P(\text{Coin})$
 - $P(\text{head or tail} \mid \text{Coin} = 0)$ and $P(\text{head or tail} \mid \text{Coin} = 1)$
 - $P(\text{weight} \mid \text{Coin} = 0)$ and $P(\text{weight} \mid \text{Coin} = 1)$

CSE 5523: Naïve Bayes



THE OHIO STATE UNIVERSITY

HW & Others

- HW-2 is due on 10/1.
- Have you been reading the textbook?

Generative vs. Discriminative learning

- Training data: $D_{tr} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$
- Goal: fit $p(X, Y)$ to D_{tr}
 - $p(Y)p(\mathbf{X}|Y)$: generative learning
 - $p(\mathbf{X})p(Y|\mathbf{X})$: discriminative learning
 - Usually not estimating $p(\mathbf{X})$
- Prediction
 - Given a test example \mathbf{x}
 - $\hat{y} = \max_c p(Y = c|\mathbf{x})$; the Bayes Optimal Classifier

Estimating “Parameters of Probability” from Data

- Maximum likelihood estimation (MLE)

- $\widehat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\mathcal{D}; \boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N); \boldsymbol{\theta}) =$
 $\underset{\boldsymbol{\theta}}{\operatorname{argmax}} p((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N); \boldsymbol{\theta}) \times \dots p((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N); \boldsymbol{\theta}) =$
 $\underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_n p((\mathbf{x}_n, y_n); \boldsymbol{\theta})$

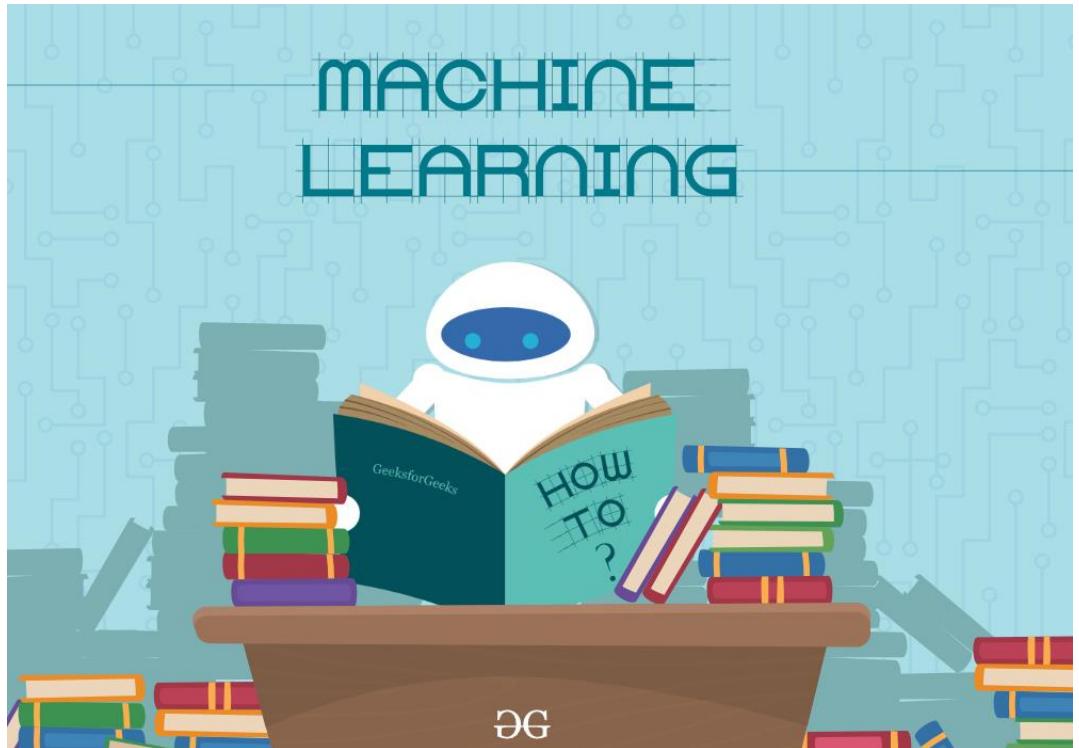
- Maximum a posteriori estimation (MAP)

- $\widehat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} P(\boldsymbol{\theta} | \mathcal{D}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\boldsymbol{\theta}) \times p((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) | \boldsymbol{\theta}) =$
 $\underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\boldsymbol{\theta}) \times \prod_n p((\mathbf{x}_n, y_n) | \boldsymbol{\theta})$

Today

Naïve Bayes

- Review
- Inference
- Advanced study



Naïve Bayes

- **Training data:** $D_{tr} = \{ (x_i \in \mathbb{R}^D, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$
 - Each dimension of x_i can be binary, categorical, count, or numerical (real)
- **Goal:** construct $p(Y = c|x)$ for $\hat{y} = \max_c p(Y = c|x)$

Naïve Bayes

- **Training data:** $D_{tr} = \{ (x_i \in \mathbb{R}^D, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$
 - Each dimension of x_i can be binary, categorical, count, or numerical (real)
- **Goal:** construct $p(Y = c|x)$ for $\hat{y} = \max_c p(Y = c|x)$
- **How? Bayes' rules**
 - $p(Y = c|x) = \frac{P(x|Y=c)p(Y=c)}{P(x)} \propto p(x|Y=c)p(Y=c)$
 - $\operatorname{argmax}_c p(Y = c|x) = \operatorname{argmax}_c p(x|Y=c)p(Y=c)$
- **Key:** what probability models for $p(x|Y = c; \theta)$ and $p(Y = c; \theta)$?

Naïve Bayes: modeling

- $p(Y = c; \theta)$: Bernoulli or categorical
- $p(x|Y = c; \theta)$: **conditional independence assumption:**
 - $p(\textcolor{red}{x}|Y = c) = p(x[1], x[2], x[3], \dots, x[D] | Y = c)$
 - $= p(x[1] | Y = c)p(x[2] | Y = c) \dots \dots p(x[D] | Y = c) = \prod_{d=1}^D p(x[d] | Y = c)$

Naïve Bayes: how to estimate parameters?

x is 3-dimensional car profile; $D_{tr} = \{ (x_i, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$

| |
|-------------|
| Maintained? |
| Color? |
| Mileage? |

$$\begin{aligned}(x|Y = c) \\ &= p(x[1] | Y = c) \\ &\times p(x[2] | Y = c) \\ &\times p(x[3] | Y = c)\end{aligned}$$

MLE or MAP?

Reading

- See “11. Naive_Bayes (additional reading).PPT” for the comparison on MLE and MAP for Naïve Bayes

Naïve Bayes: formal learning

- MLE:
 - Maximize the “joint probability” of the **IID** data instances

$$\log p(\{x_1, y_1\}, \dots) = \log \prod_{i=1}^N p(x_i, y_i) = \sum_{i=1}^N \log p(x_i, y_i) = \sum_{i=1}^N \log p(x_i|y_i)p(y_i)$$

Naïve Bayes: learning

- MLE:

- Maximize the “joint probability” of the **IID** data instances
- Separately optimize each probability model: each has independent parameters and can have different models!

$$\sum_{d=1}^D \sum_{i=1}^N \log p(\mathbf{x}_i[d] | y_i) + \sum_{i=1}^N \log p(y_i)$$
$$= \boxed{\sum_{i=1}^N \log p(\mathbf{x}_i[1] | y_i)} + \dots \dots + \boxed{\sum_{i=1}^N \log p(\mathbf{x}_i[D] | y_i)} + \boxed{\sum_{i=1}^N \log p(y_i)}$$

Categorical, Bernoulli, or 1-D Gaussian for each class c

Categorical or Bernoulli

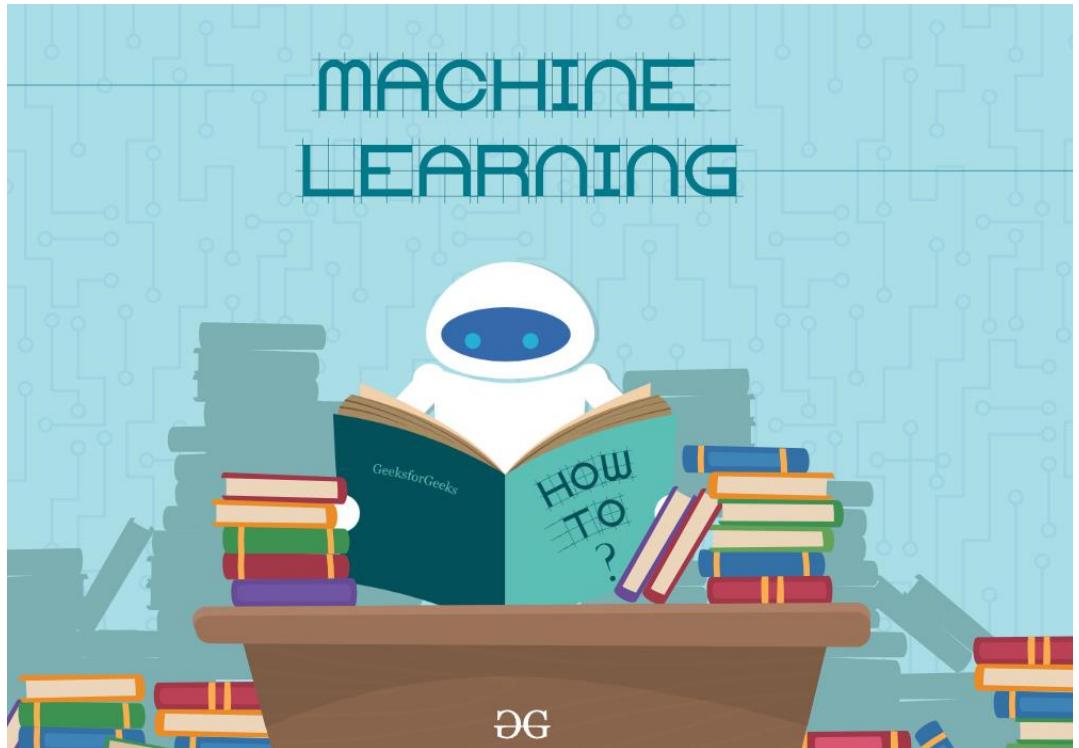
Naïve Bayes: learning

- MLE: for each conditional term $\sum_{i=1}^N \log p(\mathbf{x}_i[d] | y_i)$, we must build a conditional probability model for each class c
 - $p(\mathbf{x}[d] | Y = 1), p(\mathbf{x}[d] | Y = 2), \dots \dots, p(\mathbf{x}[d] | Y = C)$
- How?
 - Look at data instances of class c
 - Pick one feature dimension d
 - Perform MLE over $\mathbf{x}_i[d]$ whose $y_i = c$ (i.e., $\{(\mathbf{x}_i, y_i = c)\}$) to build $p(x[d] | Y = c)$
 - If $x[d]$ is continuous (e.g., weight), use a Gaussian
 - If $x[d]$ is binary, use Bernoulli; if $x[d]$ is categorical, use categorical
 - Of course, other distributions are applicable!

Today

Naïve Bayes

- Review
- Inference
- Advanced study



Naïve Bayes: inference

- Classification rule for a given x :

$$\begin{aligned} \underset{c}{\operatorname{argmax}} p(Y = c|x) &= \underset{c}{\operatorname{argmax}} p(x|Y = c)p(Y = c) \\ &= \underset{c}{\operatorname{argmax}} \prod_{d=1}^D p(x[d] | Y = c)p(Y = c) \\ &= \underset{c}{\operatorname{argmax}} \left\{ \sum_{d=1}^D \log p(x[d] | Y = c) \right\} + \log p(Y = c) \end{aligned}$$

Naïve Bayes: inference

- **Classification rule for a given x :**

- $$\begin{aligned} \underset{c}{\operatorname{argmax}} p(Y = c|x) &= \underset{c}{\operatorname{argmax}} p(x|Y = c)p(Y = c) \\ &= \underset{c}{\operatorname{argmax}} \prod_{d=1}^D p(x[d] | Y = c) p(Y = c) \\ &= \underset{c}{\operatorname{argmax}} \left\{ \sum_{d=1}^D \log p(x[d] | Y = c) \right\} + \log p(Y = c) \end{aligned}$$

- Implementation:

- For each c , compute $\prod_{d=1}^D p(x[d] | Y = c) p(Y = c)$
- Usually, we compute $\sum_{d=1}^D \log p(x[d] | Y = c) + \log p(Y = c)$ to prevent overflow/underflow
- Output the class c whose above value is the largest

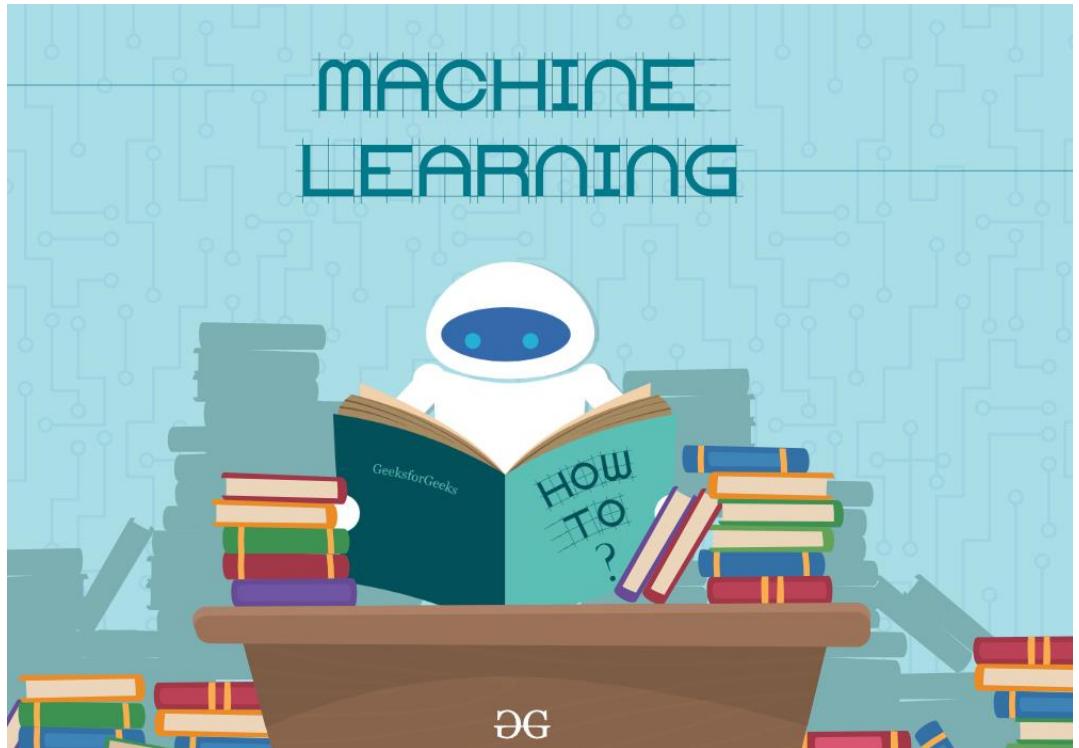
Naïve Bayes: short summary

- **Conditional independence**
 - Can effectively model random variables with continuous and discrete values
 - Conditional independence saves the number of parameters to learn
 - With CI: $P(\text{head or tail} \mid \text{Coin} = 0)$, $P(\text{heavy or light} \mid \text{Coin} = 0)$
 - Without CI: $P(\text{head or tail, heavy or light} \mid \text{Coin} = 0)$
 - Conditional independence saves the inference time
- **Bayes rules**
 - Modeling how the cause (e.g., coin type) leads to the effects (e.g., coin weight, head/tail)

Today

Naïve Bayes

- Review
- Inference
- Advanced study



Naïve Bayes: as linear classifier

- Assume that $D_{tr} = \{ (\mathbf{x}_i, y_i \in \{+1, -1\}) \}_{i=1}^N$
- $\underset{c}{\operatorname{argmax}} \prod_{d=1}^D p(x[d] | Y = c) p(Y = c)$ can be rewritten as $\operatorname{sign}(\mathbf{w}^T \mathbf{x} + b)$ if
 - $p(x[d] | Y = c)$ is Gaussian for all class c , all dimension d
 - For each dimension d , $\sigma_{d,c}^2$ is the same for all class c

Naïve Bayes: as linear classifier

- Assume that $D_{tr} = \{ (\mathbf{x}_i, y_i \in \{+1, -1\}) \}_{i=1}^N$
- $\underset{c}{\operatorname{argmax}} \prod_{d=1}^D p(x[d] | Y = c) p(Y = c)$ can be rewritten as $\operatorname{sign}(\mathbf{w}^T \mathbf{x} + b)$ if
 - $p(x[d] | Y = c)$ is Gaussian for all class c , all dimension d
 - For each dimension d , $\sigma_{d,c}^2$ is the same for all class c
- $\log p(\mathbf{x}|Y = +1)p(Y = +1) = \sum_{d=1}^D \log p(x[d] | Y = +1) + \log p(Y = +1)$

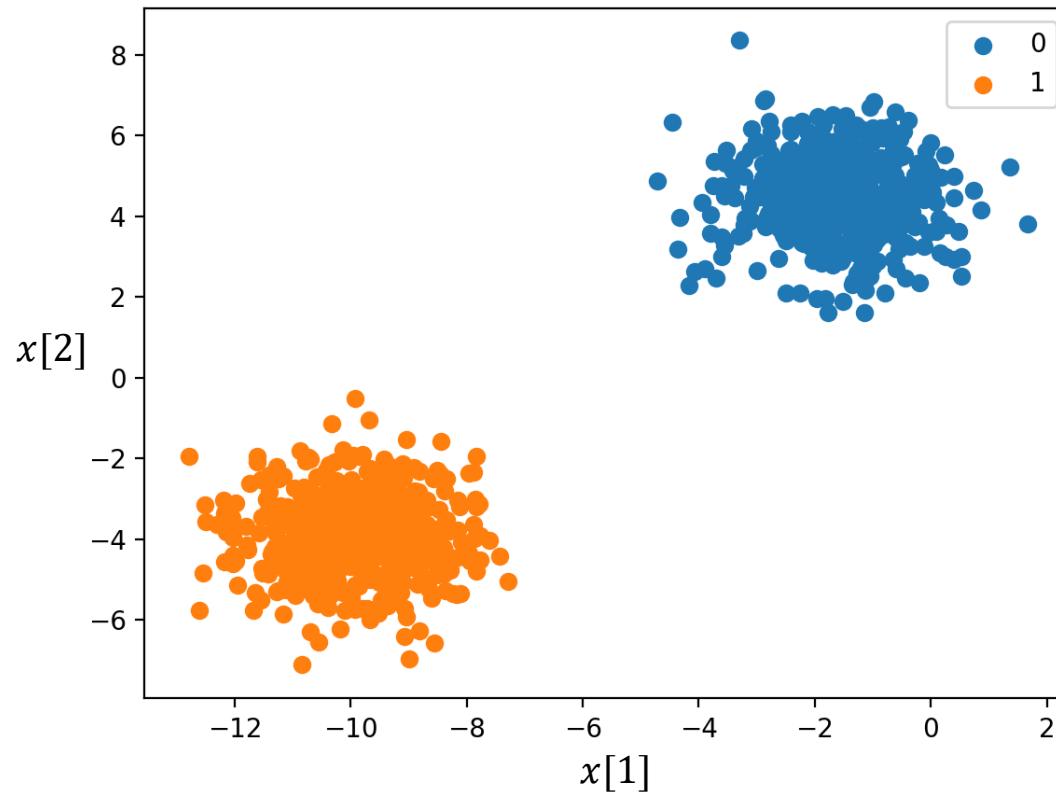
Naïve Bayes: as linear classifier

- Assume that $D_{tr} = \{ (\mathbf{x}_i, y_i \in \{+1, -1\}) \}_{i=1}^N$
- $\underset{c}{\operatorname{argmax}} \prod_{d=1}^D p(x[d] | Y = c) p(Y = c)$ can be rewritten as $\operatorname{sign}(\mathbf{w}^T \mathbf{x} + b)$ if
 - $p(x[d] | Y = c)$ is Gaussian for all class c , all dimension d
 - For each dimension d , $\sigma_{d,c}^2$ is the same for all class c
- $\log p(\mathbf{x}|Y = +1)p(Y = +1) = \sum_{d=1}^D \log p(x[d] | Y = +1) + \log p(Y = +1)$
 - $\text{Const} + \sum_{d=1}^D -\frac{(x[d] - \mu_{d,c})^2}{2\sigma_{d,c}^2} + \log \lambda = \text{Const} + \sum_{d=1}^D \frac{-x[d]^2 + 2\mu_{d,c}x[d] - \mu_{d,c}^2}{2\sigma_d^2} + \log \lambda$
- Predict $+1$ if ...

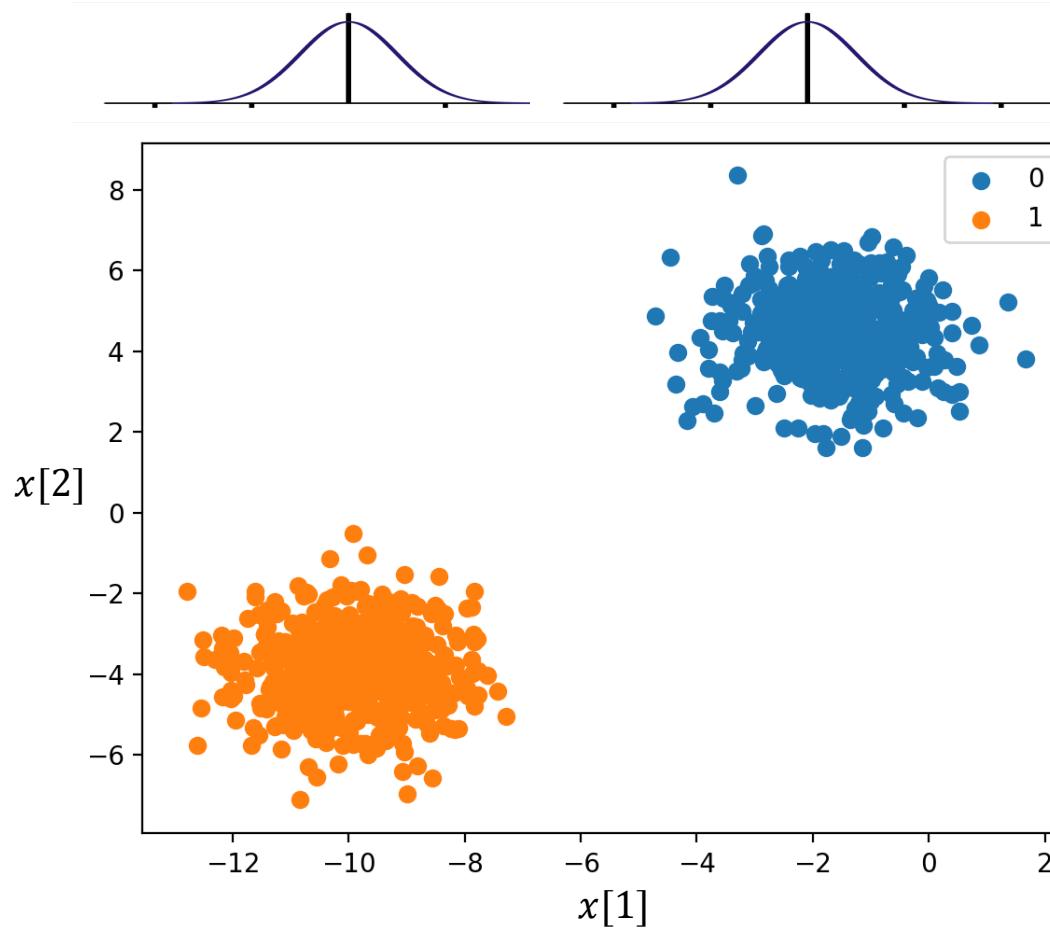
Naïve Bayes: as linear classifier

- Assume that $D_{tr} = \{ (\mathbf{x}_i, y_i \in \{+1, -1\}) \}_{i=1}^N$
- $\underset{c}{\operatorname{argmax}} \prod_{d=1}^D p(x[d] | Y = c) p(Y = c)$ can be rewritten as $\operatorname{sign}(\mathbf{w}^T \mathbf{x} + b)$ if
 - $p(x[d] | Y = c)$ is Gaussian for all class c , all dimension d
 - For each dimension d , $\sigma_{d,c}^2$ is the same for all class c
- $\log p(\mathbf{x}|Y = +1)p(Y = +1) = \sum_{d=1}^D \log p(x[d] | Y = +1) + \log p(Y = +1)$
 - $\text{Const} + \sum_{d=1}^D -\frac{(x[d] - \mu_{d,+1})^2}{2\sigma_{d,+1}^2} + \log \lambda = \text{Const} + \sum_{d=1}^D \frac{-x[d]^2 + 2\mu_{d,+1}x[d] - \mu_{d,+1}^2}{2\sigma_d^2} + \log \lambda$
- Predict $+1$ if $\sum_{d=1}^D \frac{2\mu_{d,+1}\mathbf{x}[d] - \mu_{d,+1}^2}{2\sigma_d^2} + \log \lambda > \sum_{d=1}^D \frac{2\mu_{d,-1}\mathbf{x}[d] - \mu_{d,-1}^2}{2\sigma_d^2} + \log(1 - \lambda)$

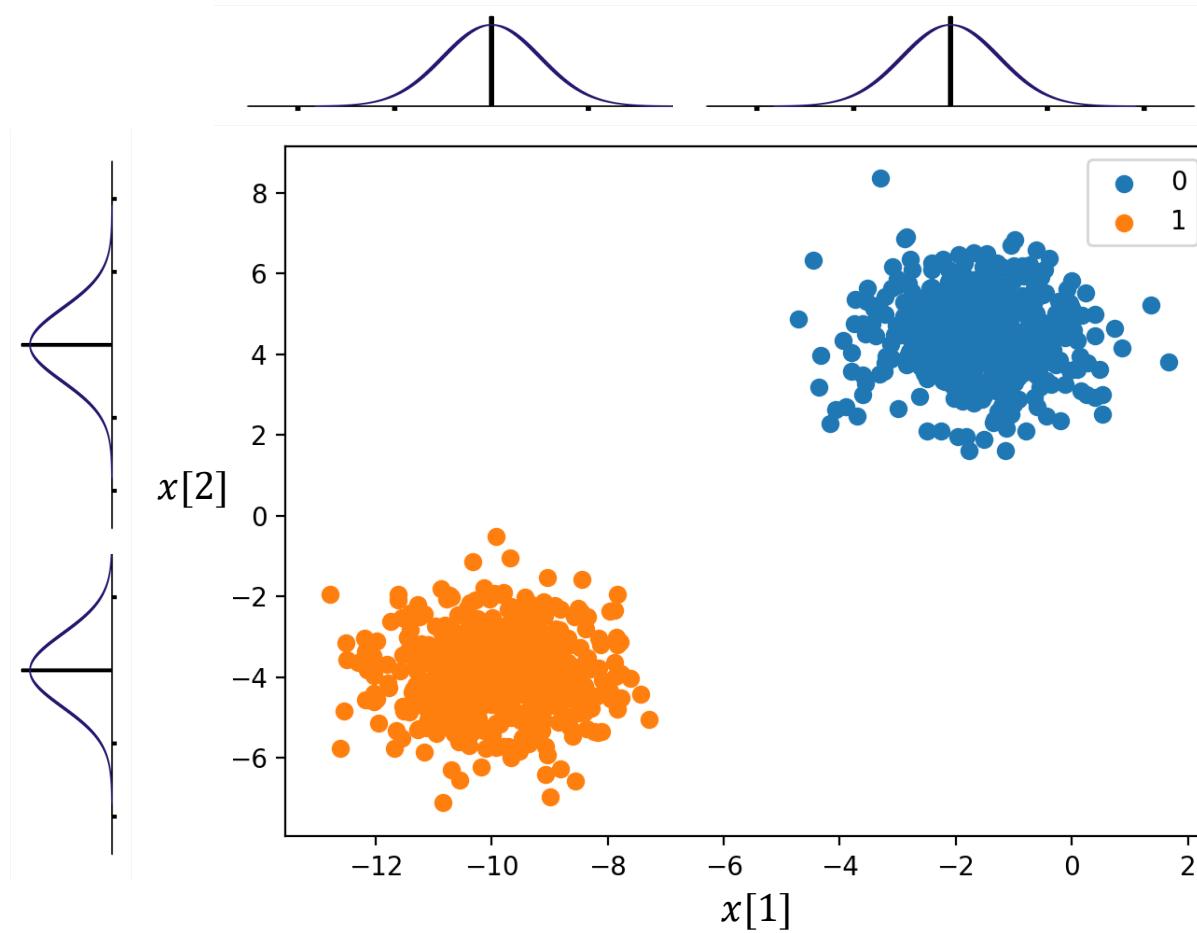
Naïve Bayes: as linear classifier



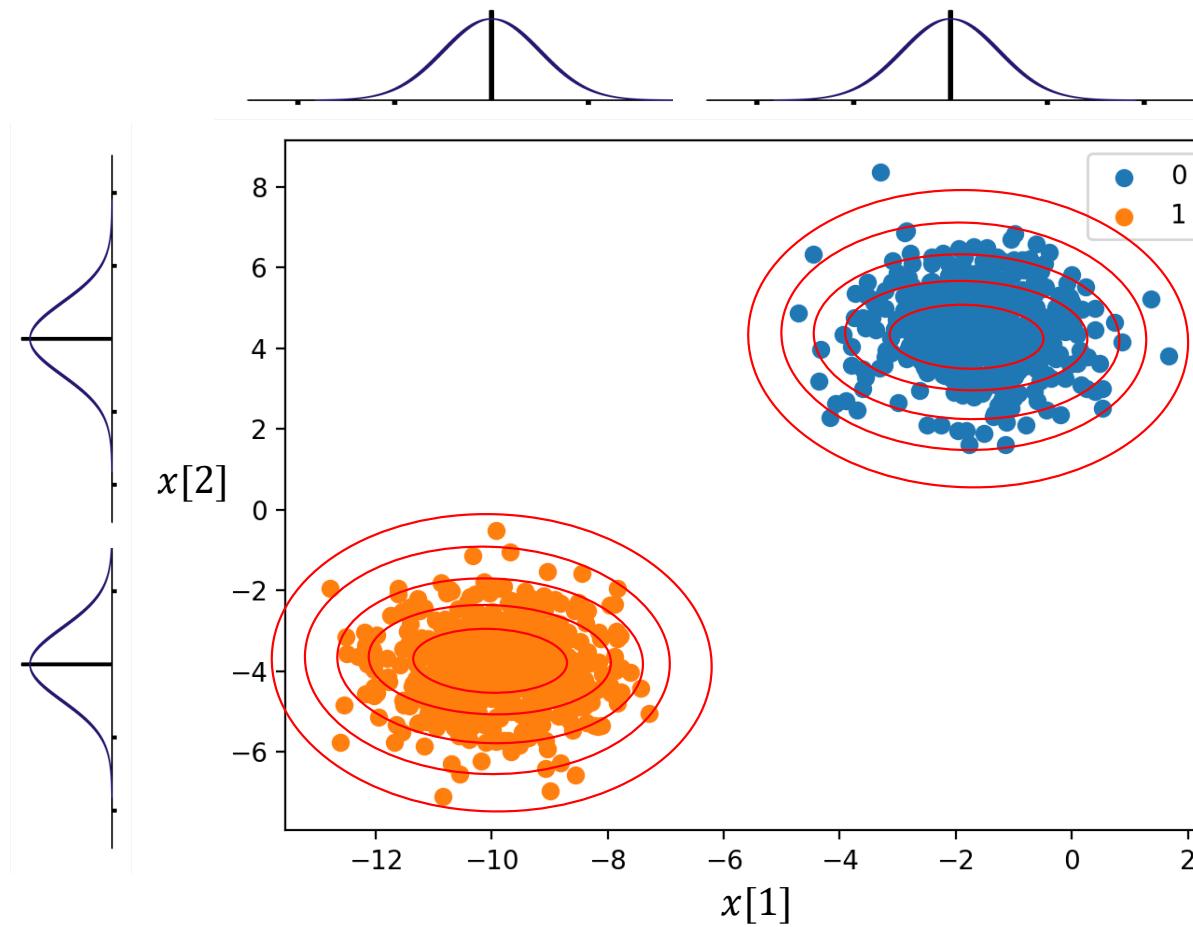
Naïve Bayes: as linear classifier



Naïve Bayes: as linear classifier



Naïve Bayes: as linear classifier



Questions?



THE OHIO STATE UNIVERSITY

Naïve Bayes: special cases for counts

- **Application:** email spam filter
 - Every email has a different number of total word counts (denoted by K)
 - Assume there is a dictionary of a fixed number of unique words D



Naïve Bayes: special cases for counts

- **Application:** email spam filter

- Every email has a different number of total word counts (denoted by K)
- Assume there is a dictionary of a fixed number of unique words D
- An email can be represented by a **D-dimensional count vector x** , ignoring word orders
 - Ex: “the dog follows the cat” has the same vector as “the cat follows the dog”
 - Ex: $x[1] = \#$ of “dog” in the email; $x[2] = \#$ of “cat”; $x[3] = \#$ of “the” ... $\rightarrow x = [1, 1, 2, \dots]^T$
 - This is called the **bag of word (BoW)** representation



Naïve Bayes: special cases for counts

- **How to represent $p(\text{email} \mid Y = c)$?**
 - The probability of seeing the count vector \mathbf{x} , i.e., $p(\mathbf{x}|Y = c)$

Naïve Bayes: special cases for counts

- **How to represent $p(\text{email} \mid Y = c)$?**
 - The probability of seeing the count vector \mathbf{x} , i.e., $p(\mathbf{x}|Y = c)$
 - The generative process:
 - Pick a class c (i.e., spam or not)
 - Pick an email length at random (i.e., a K -word email)
 - For each word k **independently**, choose a unique word (index d) from the dictionary with $p(d|Y = c)$



Dear students , I updated the programming part ,



Same dice for all words
of non-spam emails!

Naïve Bayes: special cases for counts

- **How to represent $p(\text{email} \mid Y = c)$?**

- The probability of seeing the count vector \mathbf{x} , i.e., $p(\mathbf{x}|Y = c)$
- The generative process:
 - Pick a class c (i.e., spam or not)
 - Pick an email length at random (i.e., a K -word email)
 - For each word k **independently**, choose a unique word (index d) from the dictionary with $p(d|Y = c)$
 - The probability of generating \mathbf{x} of class c is

$$p(\mathbf{x}|Y = c) = \frac{K!}{x[1]! \times \dots \times x[D]!} \prod_{d=1}^D p(d|Y = c)^{x[d]}$$

The probability to write the d-th unique word

How many times of writing it?

- Parameters: $p(d|Y = c) = \lambda_c[d]$, where $\sum_{d=1}^D \lambda_c[d] = 1$; i.e., one dice for each class

Naïve Bayes: special cases for counts

- MLE:
 - Maximize the “joint probability” of the **IID** data instances

$$\begin{aligned} \log \prod_{i=1}^N p(\mathbf{x}_i, y_i) &= \sum_{i=1}^N \log p(\mathbf{x}_i, y_i) = \sum_{i=1}^N \log p(\mathbf{x}_i|y_i)p(y_i) \\ &= \sum_{i=1}^N \log p(\mathbf{x}_i|y_i) + \sum_{i=1}^N \log p(y_i) \propto \sum_{i=1}^N \log \prod_{d=1}^D p(d|y_i)^{\mathbf{x}_i[d]} + \sum_{i=1}^N \log p(y_i) \\ &= \sum_{i=1}^N \sum_{d=1}^D \mathbf{x}_i[d] \log \lambda_{y_i}[d] + \sum_{i=1}^N \log p(y_i) \end{aligned}$$

Naïve Bayes: special cases for counts

- Parameter estimation:

- $\lambda_c[d] = \frac{\sum_{i=1}^N I[y_i=c]x_i[d] + (\alpha_c[d]-1)}{\sum_{i=1}^N I[y_i=c] \sum_{d=1}^D x_i[d] + \sum_{d=1}^D (\alpha_c[d]-1)}$
- When $\alpha_c[d] = 1$, $\lambda_c[d] = \frac{\# \text{ of times word } d \text{ appears in all emails of class } c}{\# \text{ of words in all emails of class } c}$

- Prediction:

- Given a new email with x
- $\operatorname{argmax}_c \prod_{d=1}^D p(d|Y=c)^{x[d]} \times p(Y=c)$

Naïve Bayes: as linear classifier

- Assume that $D_{tr} = \{ (\mathbf{x}_i, y_i \in \{+1, -1\}) \}_{i=1}^N$
- $\underset{c}{\operatorname{argmax}} \prod_{d=1}^D p(d|Y=c)^{x[d]} p(Y=c)$ can be rewritten as $\operatorname{sign}(\mathbf{w}^T \mathbf{x} + b)$

Questions?



THE OHIO STATE UNIVERSITY

Summary: Naïve Bayes

- To apply generative learning, it is common to start with
 - How to model $p(\text{a data instance} \mid Y = y)$
 - How to “sample” or “generate” a data instance given $Y = y$
 - Identify what parameters to learn
 - Write down the joint probability
 - Perform MLE or MAP

Summary: Naïve Bayes

- **Training data:** $D_{tr} = \{ (x_i \in \mathbb{R}^D, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$
 - Each dimension of x_i can be binary, categorical, count, or numerical (real)
- **How to construct $p(Y = c|x)$? Bayes' rules**
 - $p(Y = c|x) = \frac{P(x|Y=c)p(Y=c)}{P(x)} \propto p(x|Y=c)p(Y=c)$
 - $\operatorname{argmax}_c p(Y = c|x) = \operatorname{argmax}_c p(x|Y=c)p(Y=c)$

Summary: Naïve Bayes

- $p(Y = c; \theta)$: Bernoulli or categorical
- $p(x|Y = c; \theta)$: **conditional independence assumption:**
 - $p(x|Y = c) = p(x[1], x[2], x[3], \dots, x[D] | Y = c)$
 - $= p(x[1] | Y = c)p(x[2] | Y = c) \dots \dots p(x[D] | Y = c) = \prod_{d=1}^D p(x[d] | Y = c)$

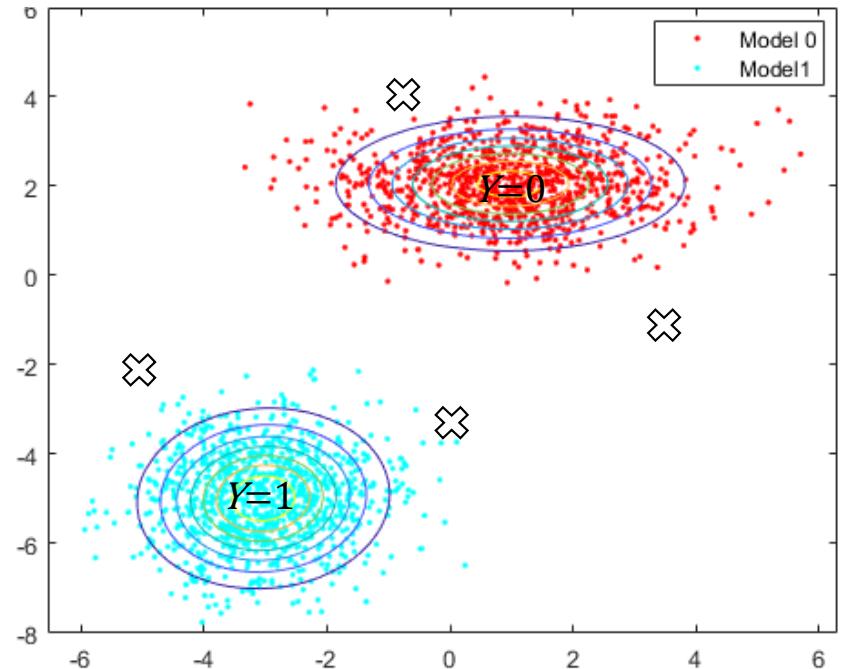
Summary: Naïve Bayes

- MLE:
 - Maximize the “joint probability” of the **IID** data instances

$$\begin{aligned} \log \prod_{i=1}^N p(\mathbf{x}_i, y_i) &= \sum_{i=1}^N \log p(\mathbf{x}_i, y_i) = \sum_{i=1}^N \log p(\mathbf{x}_i | y_i) p(y_i) \\ &= \sum_{i=1}^N \log p(\mathbf{x}_i | y_i) + \sum_{i=1}^N \log p(y_i) = \sum_{i=1}^N \log \prod_{d=1}^D p(\mathbf{x}_i[d] | y_i) + \sum_{i=1}^N \log p(y_i) \\ &= \sum_{d=1}^D \sum_{i=1}^N \log p(\mathbf{x}_i[d] | y_i) + \sum_{i=1}^N \log p(y_i) \end{aligned}$$

Summary: generative learning for classification

- Given $\{(x_i, y_i)\}_{i=1}^N$
- Assume each instance is generated by
 - $y_i \sim p(Y)$
 - $x_i \sim p(X | Y=y_i)$
- Estimate parameters (MLE)
 - $p(Y), p(X | Y=c)$ for each class c
 - Note: $p(Y) \times p(X | Y) = p(X, Y)$
- Prediction for x
 - $\underset{c}{\operatorname{argmax}} p(Y = c | x)$
 - $\underset{c}{\operatorname{argmax}} p(x | Y = c)p(Y = c)$
 - Which class is more likely to generate x ?



Naïve Bayes

x is 3-dimensional car profile; $D_{tr} = \{ (x_i, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$

Registered?

Color?

Mileage?

Naïve Bayes: MLE

x is 3-dimensional car profile; $D_{tr} = \{ (x_i, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$

| |
|-------------|
| Registered? |
| Color? |
| Mileage? |



Bernoulli distribution: $p(x[d] = 1 | Y = c) = \lambda_{d,c}$

$$\lambda_{d,c} = \frac{\sum_{i=1}^N \mathbf{1}[y_i = c] \mathbf{1}[x_i[d] = 1]}{\sum_{i=1}^N \mathbf{1}[y_i = c]}$$

Each pair of (dimension, class) has a specific coin

The generative process:

1. Sample one class c
2. Flip the corresponding coin with $\lambda_{d,c}$

Naïve Bayes: MAP

x is 3-dimensional car profile; $D_{tr} = \{ (x_i, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$

| |
|-------------|
| Registered? |
| Color? |
| Mileage? |



Bernoulli distribution: $p(x[d] = 1 | Y = c) = \lambda_{d,c}$

$$\lambda_{d,c} = \frac{\sum_{i=1}^N \mathbf{1}[y_i = c] \mathbf{1}[x_i[d] = 1] + (\alpha_{d,c} - 1)}{\sum_{i=1}^N \mathbf{1}[y_i = c] + (\alpha_{d,c} + \beta_{d,c} - 2)}$$

Each pair of (dimension, class) has a specific coin

The generative process:

1. Sample one class c
2. Flip the corresponding coin with $\lambda_{d,c}$

Naïve Bayes: MLE

x is 3-dimensional car profile; $D_{tr} = \{ (x_i, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$

| |
|-------------|
| Registered? |
| Color? |
| Mileage? |



Categorical distribution: $p(x[d] = m | Y = c) = \lambda_{d,c}[m]$

$$\lambda_{d,c}[m] = \frac{\sum_{i=1}^N \mathbf{1}[y_i = c] \mathbf{1}[x_i[d] = m]}{\sum_{i=1}^N \mathbf{1}[y_i = c]}$$

category m

Each pair of (dimension, class) has a specific dice

The generative process:

1. Sample one class c
2. Throw the corresponding dice with $\lambda_{d,c}$

Naïve Bayes: MAP

x is 3-dimensional car profile; $D_{tr} = \{ (x_i, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$

| |
|-------------|
| Registered? |
| Color? |
| Mileage? |



Categorical distribution: $p(x[d] = m | Y = c) = \lambda_{d,c}[m]$

$$\lambda_{d,c}[m] = \frac{\sum_{i=1}^N \mathbf{1}[y_i = c] \mathbf{1}[x_i[d] = m] + (\alpha_{d,c}[m] - 1)}{\sum_{i=1}^N \mathbf{1}[y_i = c] + \sum_{m=1}^M (\alpha_{d,c}[m] - 1)}$$

category m

Each pair of (dimension, class) has a specific dice

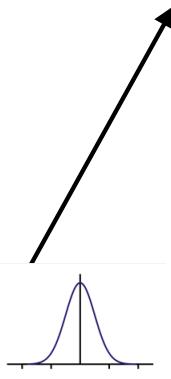
The generative process:

1. Sample one class c
2. Throw the corresponding dice with $\lambda_{d,c}$

Naïve Bayes: MLE

\mathbf{x} is 3-dimensional car profile; $D_{tr} = \{ (\mathbf{x}_i, y_i \in \{1, 2, \dots, C\}) \}_{i=1}^N$

| |
|-------------|
| Registered? |
| Color? |
| Mileage? |



Gaussian distribution: $p(x[d]|Y = c) = \mathcal{N}(\mu_{d,c}, \sigma_{d,c}^2)$

$$\mu_{d,c} = \frac{1}{N_c} \sum_{i=1}^N \mathbf{1}[y_i = c] x[d]$$

$$\sigma_{d,c}^2 = \frac{1}{N_c} \sum_{i=1}^N \mathbf{1}[y_i = c] (x[d] - \mu_{d,c})^2$$

Each pair of (dimension, class) has a specific Gaussian
The generative process:

1. Sample one class c
2. Sample from the corresponding Gaussian with $\mu_{d,c}$ and $\sigma_{d,c}^2$

Naïve Bayes: special cases for counts: MAP

- Parameter estimation:

- $\lambda_c[d] = \frac{\sum_{i=1}^N I[y_i=c]x_i[d] + (\alpha_c[d]-1)}{\sum_{i=1}^N I[y_i=c] \sum_{d=1}^D x_i[d] + \sum_{d=1}^D (\alpha_c[d]-1)}$
- When $\alpha_c[d] = 1$, $\lambda_c[d] = \frac{\text{\# of times word } d \text{ appears in all emails of class } c}{\text{\# of words in all emails of class } c}$

- Prediction:

- Given a new email with x
- $\operatorname{argmax}_c \prod_{d=1}^D p(d|Y=c)^{x[d]} \times p(Y=c)$

Naïve Bayes: special cases for counts: MLE

- Parameter estimation:

- $\lambda_c[d] = \frac{\sum_{i=1}^N I[y_i=c]x_i[d]}{\sum_{i=1}^N I[y_i=c] \sum_{d=1}^D x_i[d]}$

- $\lambda_c[d] = \frac{\text{\# of times word } d \text{ appears in all emails of class } c}{\text{\# of words in all emails of class } c}$

- Prediction:

- Given a new email with x

- $\underset{c}{\operatorname{argmax}} \prod_{d=1}^D p(d|Y=c)^{x[d]} \times p(Y=c)$

CSE 5523: Logistic Regression



THE OHIO STATE UNIVERSITY

Quick Notes

- HW2 is due next Tuesday at midnight.
- HW3 will be released by next Wednesday at midnight.
- Midterm is on 10/17 in class: you can have up to 5 papers (10 pages) of cheat sheets (A4/letter size).

Generative vs. Discriminative learning

- Training data: $D_{tr} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$
- Goal: fit $p(X, Y)$ to D_{tr}
 - $p(Y)p(\mathbf{X}|Y)$: generative learning
 - $p(\mathbf{X})p(Y|\mathbf{X})$: discriminative learning
 - Usually not estimating $p(\mathbf{X})$
- Prediction (in the MLE or MAP scenario, mostly we do)
 - Given a test example \mathbf{x}
 - $\hat{y} = \max_c p(Y = c|\mathbf{x})$; the Bayes Optimal Classifier

Naïve Bayes: generative learning

- $p(Y = c; \theta)$: Bernoulli or categorical
- $p(x|Y = c; \theta)$: **conditional independence assumption:**
 - $p(x|Y = c) = p(x[1], x[2], x[3], \dots, x[D] | Y = c)$
 - $= p(x[1] | Y = c)p(x[2] | Y = c) \dots \dots p(x[D] | Y = c) = \prod_{d=1}^D p(x[d] | Y = c)$

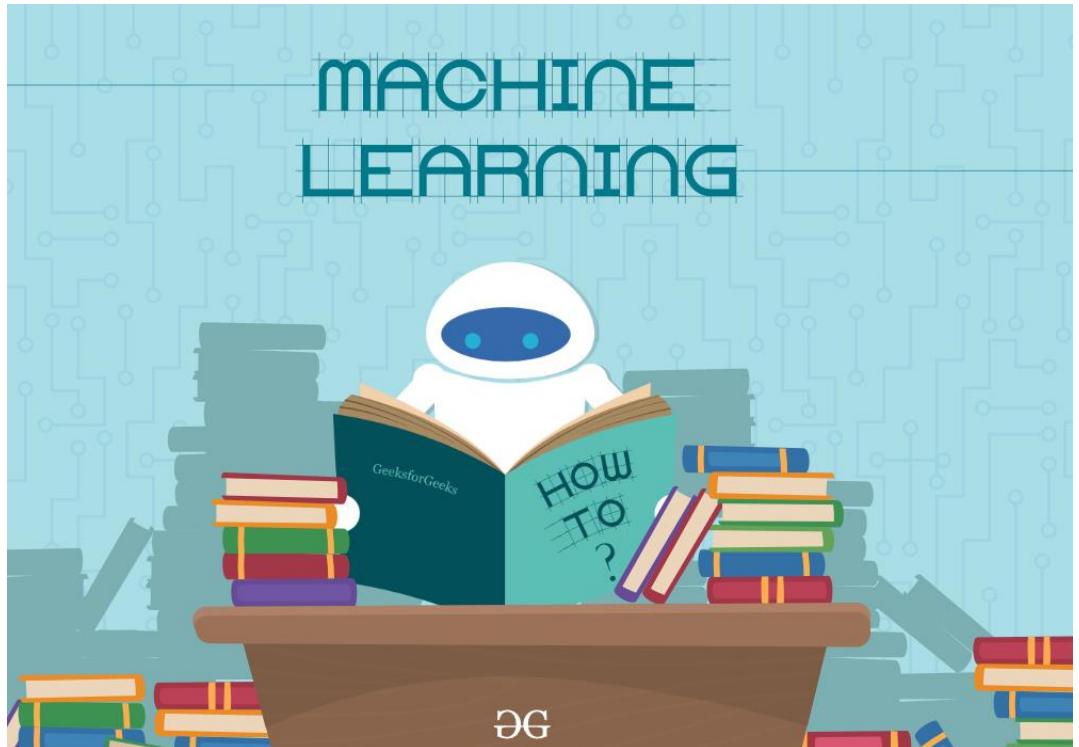
Naïve Bayes: generative learning

- $p(Y = c; \theta)$: Bernoulli or categorical
- $p(x|Y = c; \theta)$: **conditional independence assumption:**
 - $p(x|Y = c) = p(x[1], x[2], x[3], \dots, x[D] | Y = c)$
 - $= p(x[1] | Y = c)p(x[2] | Y = c) \dots \dots p(x[D] | Y = c) = \prod_{d=1}^D p(x[d] | Y = c)$
- How to learn from data?
 - So far, the distributions we use give “closed-form” solutions
 - For some other distributions, gradient descents may be needed
 - Note: For a Bernoulli distribution, the outcomes are not necessarily $\{1, 0\}$.

Today

Logistic regression

- Modeling
- Learning
- Learning (implementation)
- GDA vs. NB vs. Logistic regression
- Another perspective



Logistic regression

- **Training data:** $D_{tr} = \{ (x_i \in \mathbb{R}^D, y_i \in \{+1, -1\}) \}_{i=1}^N$
- **How to construct $p(Y = +1|x)$?**
 - Conditional Bernoulli; i.e., $p(Y = +1|x) = \lambda(x)$
 - Assumption: $\lambda(x) = \frac{1}{1+e^{-(w^T x + b)}} = \frac{1}{1+e^{-(\tilde{w}^T \tilde{x})}}$
 - Reason 1: from Naïve Bayes (see HW #3)
- **Linear or nonlinear classifier?**

Logistic regression

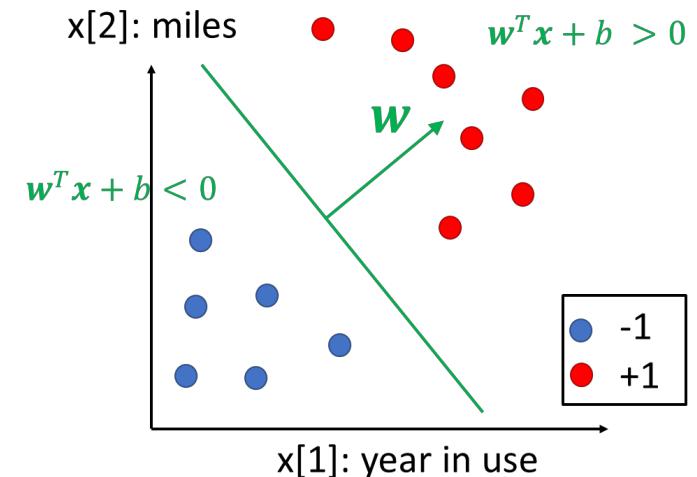
- **Training data:** $D_{tr} = \{ (x_i \in \mathbb{R}^D, y_i \in \{+1, -1\}) \}_{i=1}^N$

- **How to construct $p(Y = +1|x)$?**

- Conditional Bernoulli; i.e., $p(Y = +1|x) = \lambda(x)$
- Assumption: $\lambda(x) = \frac{1}{1+e^{-(w^T x + b)}} = \frac{1}{1+e^{-(\tilde{w}^T \tilde{x})}}$
- Reason 1: from Naïve Bayes (see HW #3)

- **Linear or nonlinear classifier?**

- Dependent on the decision boundary!
- Bayes optimal classifier predicts $+1$ when $\lambda(x) > 0.5$; i.e., when $\tilde{w}^T \tilde{x} > 0$



Logistic regression

- Reason 2: ratio for classification
 - Linear: $\tilde{\mathbf{w}}^T \tilde{\mathbf{x}} > 0$ for class +1
 - Probability: $r = \frac{p(Y = +1|\mathbf{x})}{1-p(Y = +1|\mathbf{x})} > 1$ or $\log r = \log \frac{p(Y = +1|\mathbf{x})}{1-p(Y = +1|\mathbf{x})} > 0$

Logistic regression

- Reason 2: ratio for classification
 - Linear: $\tilde{\mathbf{w}}^T \tilde{\mathbf{x}} > 0$ for class +1
 - Probability: $r = \frac{p(Y = +1|\mathbf{x})}{1-p(Y = +1|\mathbf{x})} > 1$ or $\log r = \log \frac{p(Y = +1|\mathbf{x})}{1-p(Y = +1|\mathbf{x})} > 0$
- Let $\tilde{\mathbf{w}}^T \tilde{\mathbf{x}} = \log \frac{p(Y = +1|\mathbf{x})}{1-p(Y = +1|\mathbf{x})}$
 - $e^{-(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})} = \frac{1-p(Y = +1|\mathbf{x})}{p(Y = +1|\mathbf{x})}$ and then $e^{-(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})} p(Y = +1|\mathbf{x}) = 1 - p(Y = +1|\mathbf{x})$
 - $p(Y = +1|\mathbf{x}) \left(1 + e^{-(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})}\right) = 1$
 - $p(Y = +1|\mathbf{x}) = \frac{1}{1+e^{-(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})}}$

Sigmoid function

- $p(Y = +1|\mathbf{x}) = \frac{1}{1+e^{-(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})}} = \rho(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})$

- $p(Y = -1|\mathbf{x}) = 1 - p(Y = +1|\mathbf{x})$

- $1 - \frac{1}{1+e^{-(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})}} = \frac{e^{-(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})}}{1+e^{-(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})}} = \frac{1}{1+e^{(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})}}$

- $p(\mathbf{y}|\mathbf{x}) = \frac{1}{1+e^{-\mathbf{y}(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})}} = \rho(\mathbf{y}\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})$

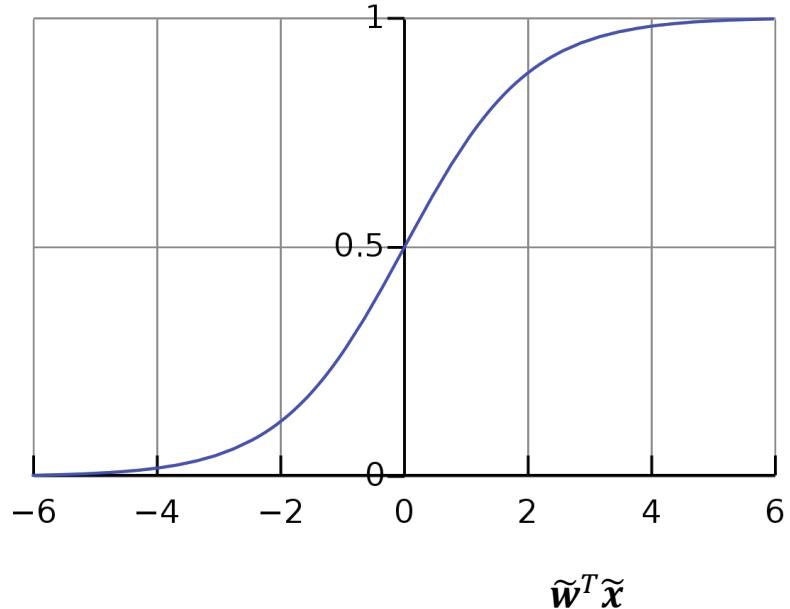
- Sigmoid function has

- Domain: \mathbb{R}

- Range: $[0, 1]$ (popular for probability)

- Monotonically increasing

Sigmoid function
 $\rho(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}})$



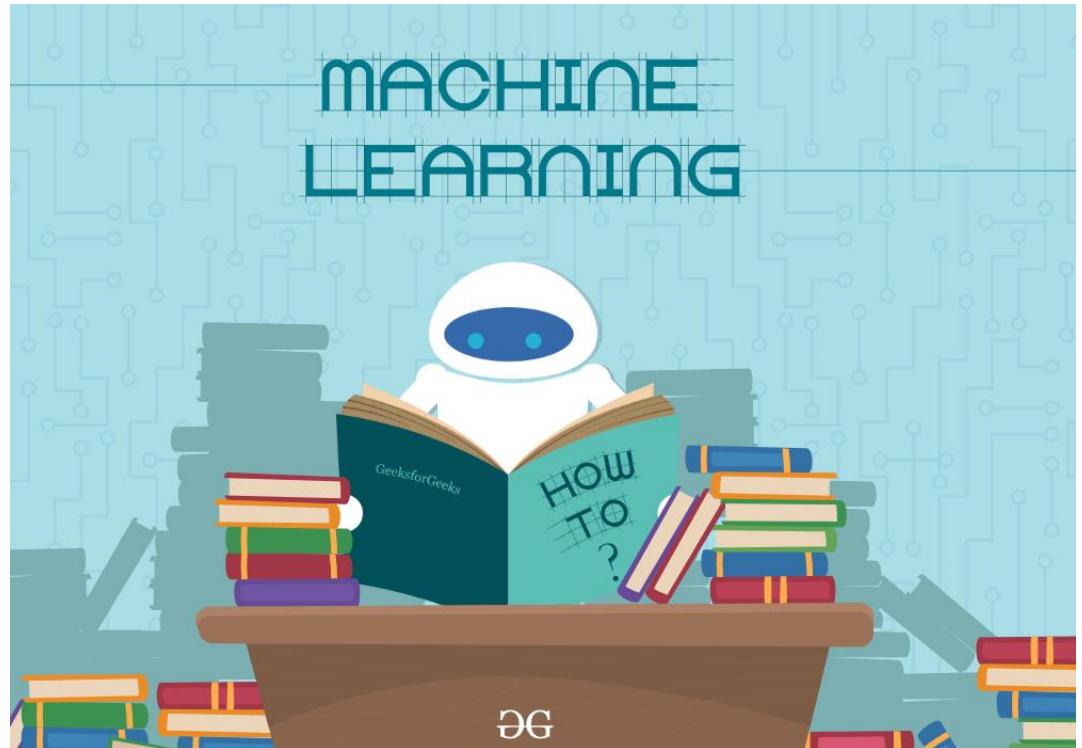
Logistic regression

- Modeling $p(y|x; \theta)$, so discriminative learning
- MLE: $\widehat{\theta} = \underset{\theta}{\operatorname{argmax}} p(D_{tr}; \theta) = \underset{\theta}{\operatorname{argmax}} \prod_n p(y_n|x_n; \theta)p(x_n)$
- MAP: $\widehat{\theta} = \underset{\theta}{\operatorname{argmax}} p(D_{tr}|\theta)p(\theta) = \underset{\theta}{\operatorname{argmax}} p(\theta) \prod_n p(y_n|x_n, \theta)p(x_n)$

Today

Logistic regression

- Modeling
- Learning
- Learning (implementation)
- GDA vs. NB vs. Logistic regression
- Another perspective



Logistic regression (MLE)

- MLE: $\hat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(D_{tr}; \boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \prod_n p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) p(\mathbf{x}_n)$

$$\begin{aligned} \log \prod_{i=1}^N p(y_i | \mathbf{x}_i; \boldsymbol{\theta}) &= \sum_{i=1}^N \log p(y_i | \mathbf{x}_i; \boldsymbol{\theta}) = \sum_{i=1}^N \log \rho(y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i) \\ &= \sum_{i=1}^N \log \frac{1}{1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i}} = \sum_{i=1}^N -\log(1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i}) \end{aligned}$$

- $\tilde{\mathbf{w}}^{\text{MLE}} = \operatorname{argmax}_{\tilde{\mathbf{w}}} \sum_{i=1}^N -\log(1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i}) = \operatorname{argmin}_{\tilde{\mathbf{w}}} \sum_{i=1}^N \log(1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i})$

○ How to solve?

Logistic regression (MLE)

- $\tilde{\mathbf{w}}^{\text{MLE}} = \underset{\tilde{\mathbf{w}}}{\operatorname{argmax}} \sum_{i=1}^N -\log \left(1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i} \right) = \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_{i=1}^N \log \left(1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i} \right)$
 - How to solve?
 - No closed-form solution for $\nabla_{\tilde{\mathbf{w}}} \left\{ \sum_{i=1}^N \log \left(1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i} \right) \right\} = 0$
 - Gradient descent!

Logistic regression (MAP)

- MAP: $\widehat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(D_{tr}; \boldsymbol{\theta}) p(\boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(\boldsymbol{\theta}) \prod_n p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) p(\mathbf{x}_n)$
- $p(\boldsymbol{\theta}) = \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$

$$\log p(\boldsymbol{\theta}) \prod_{i=1}^N p(y_i | \mathbf{x}_i; \boldsymbol{\theta}) \propto \sum_{i=1}^N -\log(1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i}) - \frac{1}{2\sigma^2} \|\tilde{\mathbf{w}}\|_2^2$$

- $\tilde{\mathbf{w}}^{\text{MAP}} = \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_{i=1}^N \log(1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i}) + \frac{1}{2\sigma^2} \|\tilde{\mathbf{w}}\|_2^2$
 - Gradient descent!

Logistic regression (MLE)

- Some textbooks use $y_i \in \{+1, 0\}$ rather than $y_i \in \{+1, -1\}$

$$\log \prod_{n=1}^N p(y_n | \boldsymbol{x}_n; \boldsymbol{\theta}) = \sum_{i=1}^N \log p(y_n | \boldsymbol{x}_n; \boldsymbol{\theta})$$

Logistic regression (MLE)

- Some textbooks use $y_i \in \{+1, 0\}$ rather than $y_i \in \{+1, -1\}$

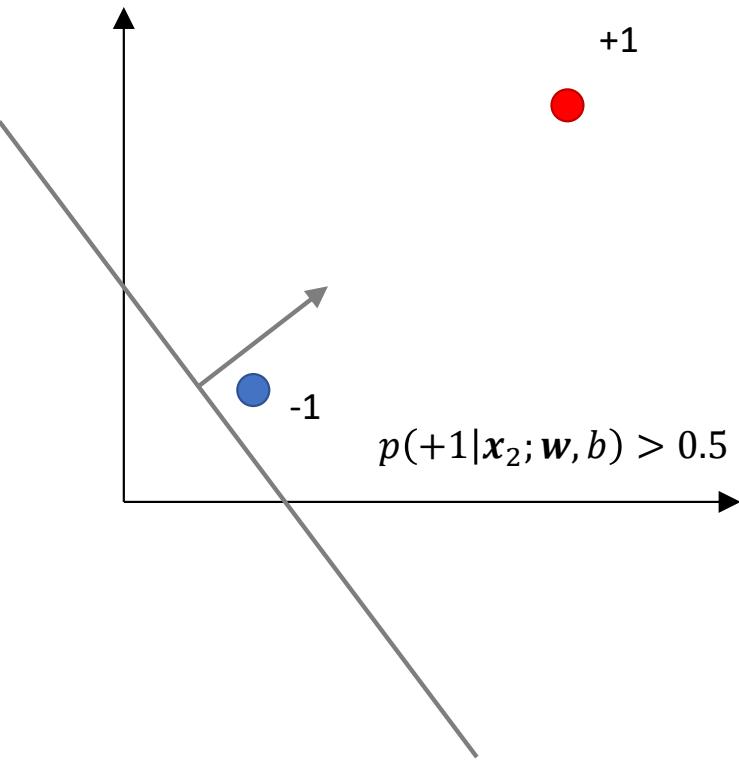
$$\begin{aligned} \log \prod_{n=1}^N p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) &= \sum_{i=1}^N \log p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) \\ &= \sum_{n=1}^N \log \left(p(+1 | \mathbf{x}_n; \boldsymbol{\theta})^{y_n} \times (1 - p(+1 | \mathbf{x}_n; \boldsymbol{\theta}))^{1-y_n} \right) \\ &= \sum_{n=1}^N y_n \log p(+1 | \mathbf{x}_n; \boldsymbol{\theta}) + (1 - y_n) \log(1 - p(+1 | \mathbf{x}_n; \boldsymbol{\theta})) \\ &= \sum_{n=1}^N y_n \times \log \frac{1}{1 + e^{-\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}} + (1 - y_n) \times \log \left(\frac{e^{-\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}}{1 + e^{-\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}} \right) \end{aligned}$$

Basically, no difference in algorithms, but difference in math expressions

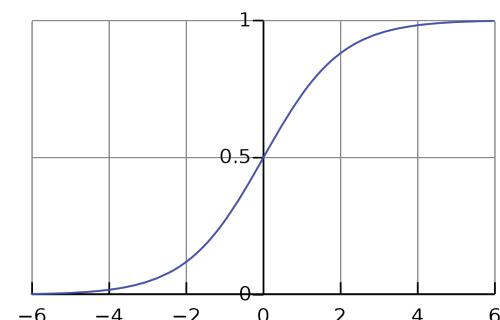
Negative of this is called
Binary cross entropy (BCE) loss

Logistic regression (demo)

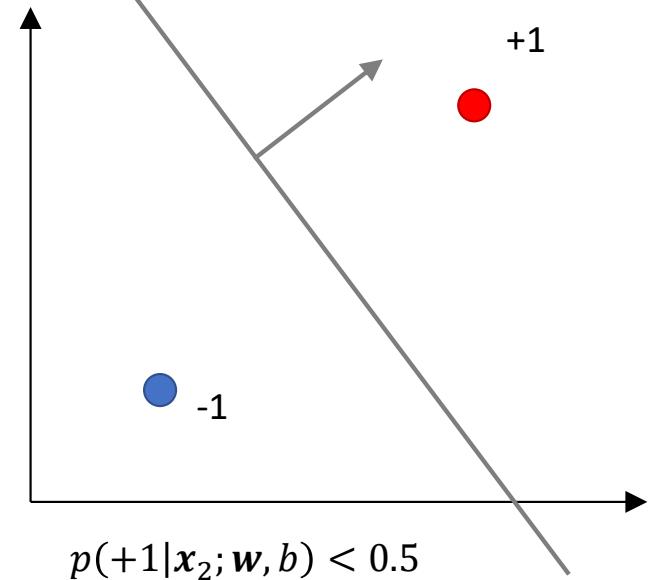
$$p(+1|x_1; \mathbf{w}, b) > 0.5$$



$$p(+1|x_2; \mathbf{w}, b) > 0.5$$

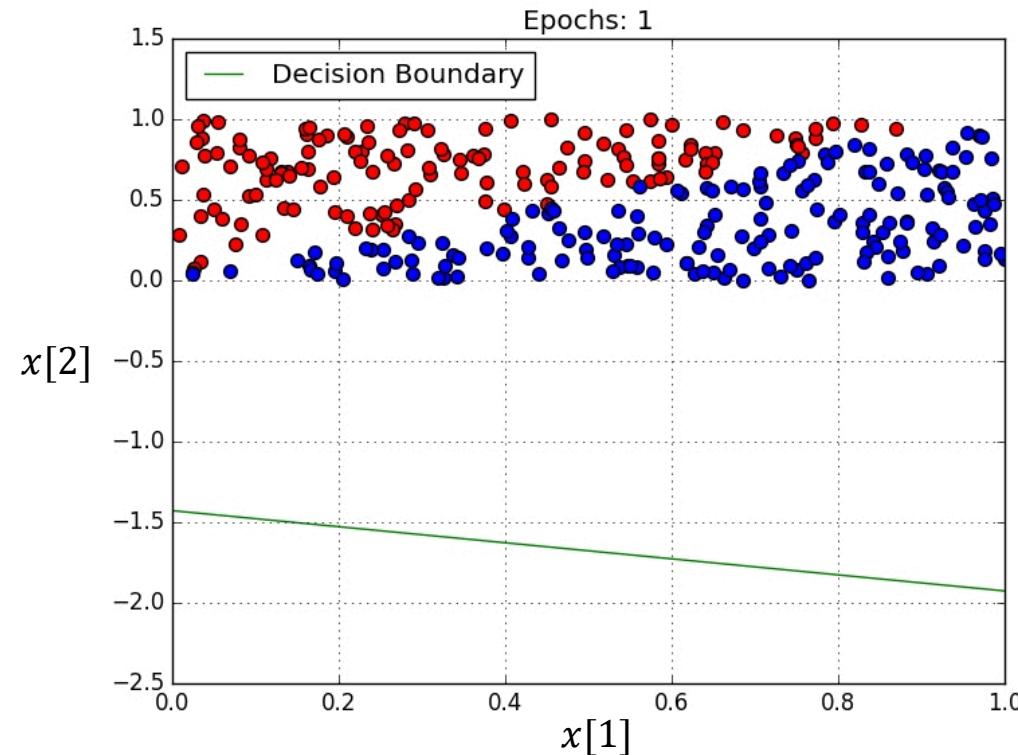


$$p(+1|x_1; \mathbf{w}, b) > 0.5$$



$$p(+1|x_2; \mathbf{w}, b) < 0.5$$

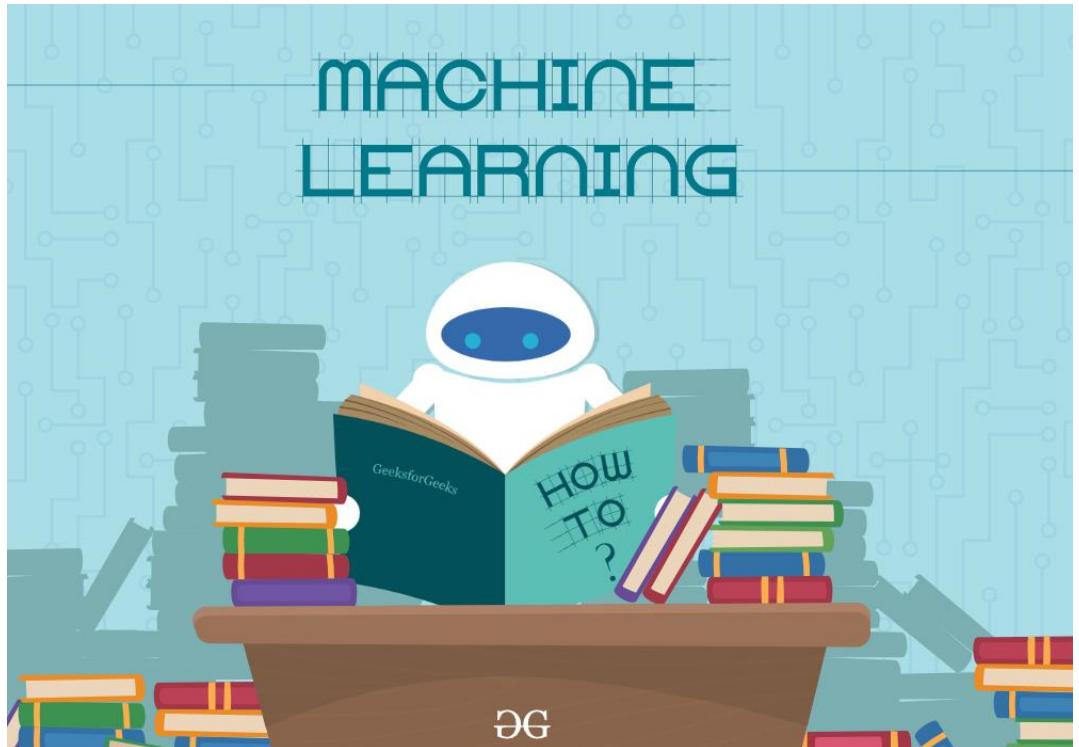
Logistic regression (demo)



Today

Logistic regression

- Modeling
- Learning
- Learning (implementation)
- GDA vs. NB vs. Logistic regression
- Another perspective



Gradient descent (GD) for logistic regression

- Initialize $\widetilde{\mathbf{w}}'$
- For $t = 1:T$
 - $\widetilde{\mathbf{w}}' \leftarrow \widetilde{\mathbf{w}}' - \eta \nabla_{\widetilde{\mathbf{w}}} E(\widetilde{\mathbf{w}}');$ $E(\widetilde{\mathbf{w}})$: the loss function you want to minimize!
 - Until T is reached or $\widetilde{\mathbf{w}}'$ converges (e.g., $\|\nabla_{\widetilde{\mathbf{w}}} E(\widetilde{\mathbf{w}}')\|_2 < \varepsilon$ or change of $E(\widetilde{\mathbf{w}}')$ is small)

Gradient descent (GD) for logistic regression

- Initialize $\widetilde{\mathbf{w}}'$
- For $t = 1:T$
 - $\widetilde{\mathbf{w}}' \leftarrow \widetilde{\mathbf{w}}' - \eta \nabla_{\widetilde{\mathbf{w}}} E(\widetilde{\mathbf{w}}');$ $E(\widetilde{\mathbf{w}})$: the loss function you want to minimize!
 - Until T is reached or $\widetilde{\mathbf{w}}'$ converges (e.g., $\|\nabla_{\widetilde{\mathbf{w}}} E(\widetilde{\mathbf{w}}')\|_2 < \varepsilon$ or change of $E(\widetilde{\mathbf{w}}')$ is small)
- Note:
 - If $E(\widetilde{\mathbf{w}}) = \frac{1}{N} \sum_i e_i(\widetilde{\mathbf{w}})$; for example, $e_i(\widetilde{\mathbf{w}})$ is the loss on the i -th example
 - $\nabla_{\widetilde{\mathbf{w}}} E(\widetilde{\mathbf{w}}) = \frac{1}{N} \sum_n \nabla_{\widetilde{\mathbf{w}}} e_n(\widetilde{\mathbf{w}})$
 - See the derivations of the gradients in later slides
 - While the gradients formulas are different for $y_i \in \{+1, -1\}$ or $y_i \in \{+1, 0\}$, the gradient values are basically the same.

Gradient descent (GD) for logistic regression

- $e_n(\tilde{\mathbf{w}}) = \log(1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n})$ for $y_i \in \{+1, -1\}$

$$\nabla_{\tilde{\mathbf{w}}} e_n(\tilde{\mathbf{w}}) = \frac{1}{1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}} \times \nabla_{\tilde{\mathbf{w}}} (1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n})$$

Derivative for Log

$$= \frac{1}{1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}} \times (-y_n \tilde{\mathbf{x}}_n) \times e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}$$

$$= \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) \times (-y_n \tilde{\mathbf{x}}_n) \times e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}$$

$$= \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) \times (-y_n \tilde{\mathbf{x}}_n) \times \frac{1 - \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)}{\rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)}$$

$$= -y_n (1 - \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)) (\tilde{\mathbf{x}}_n)$$

Gradient descent (GD) for logistic regression

- $e_n(\tilde{\mathbf{w}}) = \log(1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n})$ for $y_i \in \{+1, -1\}$

$$\nabla_{\tilde{\mathbf{w}}} e_n(\tilde{\mathbf{w}}) = \frac{1}{1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}} \times \nabla_{\tilde{\mathbf{w}}} (1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n})$$

$$= \frac{1}{1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}} \times (-y_n \tilde{\mathbf{x}}_n) \times e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}$$

$$= \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) \times (-y_n \tilde{\mathbf{x}}_n) \times e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}$$

$$= \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) \times (-y_n \tilde{\mathbf{x}}_n) \times \frac{1 - \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)}{\rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)}$$

$$= -y_n (1 - \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)) (\tilde{\mathbf{x}}_n)$$

Derivative for constants and exponential

Gradient descent (GD) for logistic regression

- $e_n(\tilde{\mathbf{w}}) = \log(1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n})$ for $y_i \in \{+1, -1\}$

$$\nabla_{\tilde{\mathbf{w}}} e_n(\tilde{\mathbf{w}}) = \frac{1}{1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}} \times \nabla_{\tilde{\mathbf{w}}} (1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n})$$

$$= \frac{1}{1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}} \times (-y_n \tilde{\mathbf{x}}_n) \times e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}$$

$$= \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) \times (-y_n \tilde{\mathbf{x}}_n) \times e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}$$

$$= \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) \times (-y_n \tilde{\mathbf{x}}_n) \times \frac{1 - \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)}{\rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)}$$

$$= -y_n (1 - \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)) (\tilde{\mathbf{x}}_n)$$

The first term is just sigmoid

Gradient descent (GD) for logistic regression

- $e_n(\tilde{\mathbf{w}}) = \log(1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n})$ for $y_i \in \{+1, -1\}$

$$\begin{aligned}\nabla_{\tilde{\mathbf{w}}} e_n(\tilde{\mathbf{w}}) &= \frac{1}{1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}} \times \nabla_{\tilde{\mathbf{w}}} (1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}) \\ &= \frac{1}{1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}} \times (-y_n \tilde{\mathbf{x}}_n) \times e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n} \\ &= \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) \times (-y_n \tilde{\mathbf{x}}_n) \times e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n} \\ &= \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) \times (-y_n \tilde{\mathbf{x}}_n) \times \frac{1 - \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)}{\rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)} \\ &= -y_n (1 - \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)) (\tilde{\mathbf{x}}_n)\end{aligned}$$

The last term can be re-written

Gradient descent (GD) for logistic regression

- $e_n(\tilde{\mathbf{w}}) = \log(1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n})$ for $y_i \in \{+1, -1\}$

$$\begin{aligned}\nabla_{\tilde{\mathbf{w}}} e_n(\tilde{\mathbf{w}}) &= \frac{1}{1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}} \times \nabla_{\tilde{\mathbf{w}}} (1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}) \\ &= \frac{1}{1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}} \times (-y_n \tilde{\mathbf{x}}_n) \times e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n} \\ &= \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) \times (-y_n \tilde{\mathbf{x}}_n) \times e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n} \\ &= \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) \times (-y_n \tilde{\mathbf{x}}_n) \times \frac{1 - \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)}{\rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)} \\ &= -y_n \boxed{(1 - \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n))(\tilde{\mathbf{x}}_n)}\end{aligned}$$

↑
Softer update

Gradient descent (GD) for logistic regression

- $e_n(\tilde{\mathbf{w}}) = \log(1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n})$ for $y_i \in \{+1, -1\}$

$$\begin{aligned}\nabla_{\tilde{\mathbf{w}}} e_n(\tilde{\mathbf{w}}) &= \frac{1}{1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}} \times \nabla_{\tilde{\mathbf{w}}} (1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}) \\ &= \frac{1}{1 + e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n}} \times (-y_n \tilde{\mathbf{x}}_n) \times e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n} \\ &= \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) \times (-y_n \tilde{\mathbf{x}}_n) \times e^{-y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n} \\ &= \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) \times (-y_n \tilde{\mathbf{x}}_n) \times \frac{1 - \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)}{\rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)} \\ &= -y_n (1 - \rho(y_n \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n)) (\tilde{\mathbf{x}}_n)\end{aligned}$$

↑
Softer update

Logistic regression:

$$\begin{aligned}\tilde{\mathbf{w}}' &\leftarrow \tilde{\mathbf{w}}' + \eta \sum_n y_n (1 - \rho(y_n \tilde{\mathbf{w}}'^T \tilde{\mathbf{x}}_n)) \tilde{\mathbf{x}}_n \\ \tilde{\mathbf{w}}' &\leftarrow \tilde{\mathbf{w}}' + \eta \sum_n y_n (1 - p(y_n | \tilde{\mathbf{x}}_n)) \tilde{\mathbf{x}}_n\end{aligned}$$

Perceptron:

$$\tilde{\mathbf{w}}' \leftarrow \tilde{\mathbf{w}}' + \eta y_n \tilde{\mathbf{x}}_n$$

Gradient descent (GD) for logistic regression

- $e_n(\tilde{\mathbf{w}}) = -[y_n \log p(+1|\mathbf{x}_n; \boldsymbol{\theta}) + (1 - y_n) \log(1 - p(+1|\mathbf{x}_n; \boldsymbol{\theta}))]$ for $y_i \in \{+1, 0\}$

$$\begin{aligned}\nabla_{\tilde{\mathbf{w}}} e_n(\tilde{\mathbf{w}}) &= -[y_n \nabla_{\tilde{\mathbf{w}}} \log p(+1|\mathbf{x}_n; \boldsymbol{\theta}) + (1 - y_n) \nabla_{\tilde{\mathbf{w}}} \log(1 - p(+1|\mathbf{x}_n; \boldsymbol{\theta}))] \\ &= -[y_n \nabla_{\tilde{\mathbf{w}}} \log \rho(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n) + (1 - y_n) \nabla_{\tilde{\mathbf{w}}} \log(1 - \rho(\tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_n))] \\ &\stackrel{=} {-} [y_n \times (1 - \rho(\mathbf{w}^T \mathbf{x}_n)) \mathbf{x}_n - (1 - y_n) \times \rho(\mathbf{w}^T \mathbf{x}_n) \mathbf{x}_n] \\ &= -(y_n - \rho(\mathbf{w}^T \mathbf{x}_n)) \mathbf{x}_n\end{aligned}$$

$$\text{“=}” : \frac{d \rho(a)}{d a} = \frac{d \left(\frac{1}{1+e^{-a}} \right)}{da} = \frac{e^{-a}}{(1+e^{-a})^2} = \frac{1}{1+e^{-a}} \times \frac{e^{-a}}{1+e^{-a}} = \rho(a)(1 - \rho(a))$$

Stochastic gradient descent (SGD) for logistic regression

- Initialize $\widetilde{\mathbf{w}}'$
- For $t = 1:T$
 - Randomly sample an example (\mathbf{x}_i, y_i) from $D_{tr} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$
 - $\widetilde{\mathbf{w}}' \leftarrow \widetilde{\mathbf{w}}' - \eta \nabla_{\widetilde{\mathbf{w}}} e_i(\widetilde{\mathbf{w}}')$
 - Until T is reached or $\widetilde{\mathbf{w}}'$ converges

Caution

- Caution:
 - The learning rate η must be chosen properly for convergence. It can be dynamic.

Short summary: Logistic regression

- MLE: $\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(D_{tr}; \boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_n p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) p(\mathbf{x}_n)$

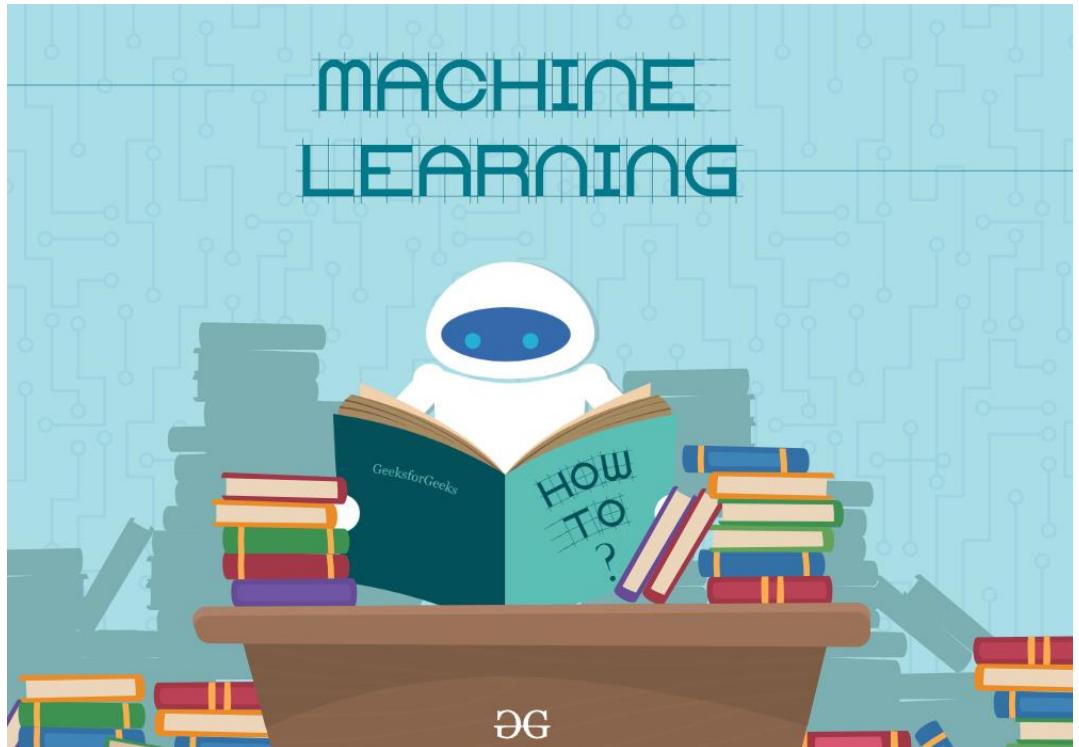
$$\begin{aligned} \log \prod_{i=1}^N p(y_i | \mathbf{x}_i; \boldsymbol{\theta}) &= \sum_{i=1}^N \log p(y_i | \mathbf{x}_i; \boldsymbol{\theta}) = \sum_{i=1}^N \log \rho(y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i) \\ &= \sum_{i=1}^N \log \frac{1}{1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i}} = \sum_{i=1}^N -\log(1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i}) \end{aligned}$$

- How to maximize?
 - Gradient descent!

Today

Logistic regression

- Modeling
- Learning
- Learning (implementation)
- GDA vs. NB vs. Logistic regression
- Another perspective



GDA vs. NB vs. Logistic regression

- Training data: $D_{tr} = \{ (\mathbf{x}_i \in \mathbb{R}^D, y_i \in \{+1, -1\}) \}_{i=1}^N$
- Under certain assumptions, all of them will lead to a linear classifier $\text{sign}(\mathbf{w}^T \mathbf{x} + b)$
- However, the learned \mathbf{w} and b are different
 - GDA and NB need NO gradient descents; they have closed-form solutions
 - GDA and NB are suitable for streaming settings

GDA vs. NB vs. Logistic regression

- Training data: $D_{tr} = \{ (\mathbf{x}_i \in \mathbb{R}^D, y_i \in \{+1, -1\}) \}_{i=1}^N$
- Under certain assumptions, all of them will lead to a linear classifier $\text{sign}(\mathbf{w}^T \mathbf{x} + b)$
- However, the learned \mathbf{w} and b are different
 - GDA and NB need NO gradient descents; they have closed-form solutions
 - GDA and NB are suitable for streaming settings
- **Recap:**
 - A learning algorithm + A hypothesis class
 - GDA, NB, and logistic regression use the same hypothesis class but different algorithms

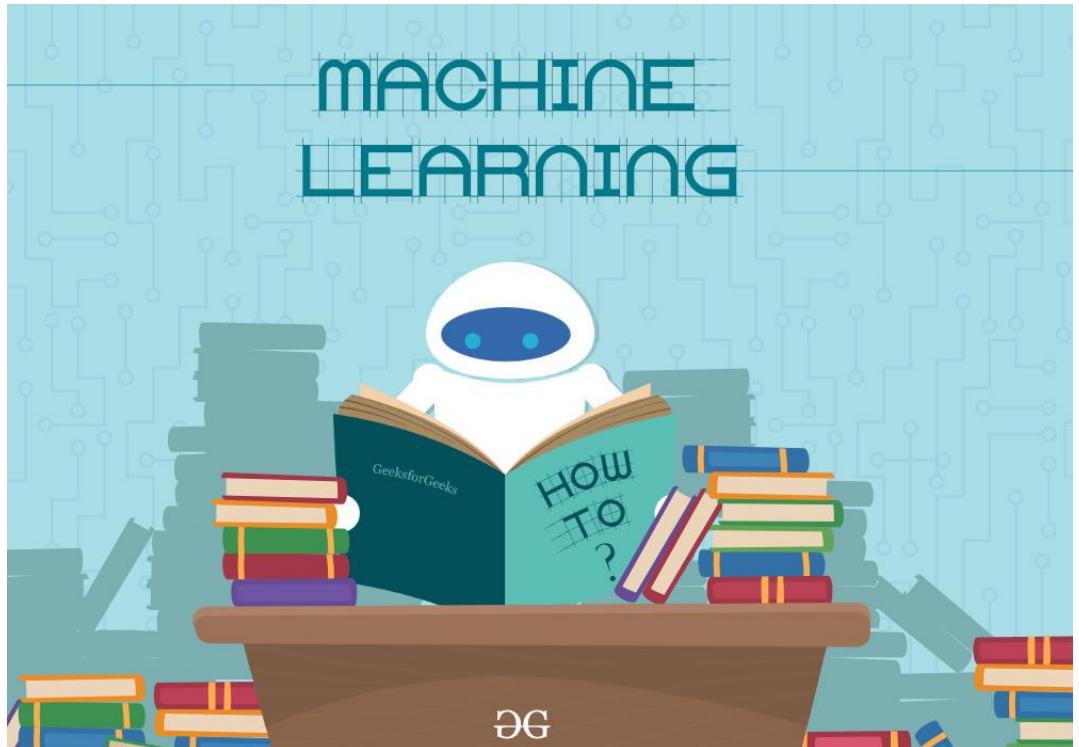
Generative vs. Discriminative learning

- Generative learning (MLE case)
 - $p(x, y; \theta) = p(x|y; \theta)p(y; \theta)$
 - Make **two modeling assumptions** (conditional data + Bernoulli or categorical)
 - If the assumption is correct, then fewer training data are needed
- Discriminative learning (MLE case)
 - $p(x, y; \theta) = p(y|x; \theta)p(x)$
 - Make **one modeling assumptions** (conditional Bernoulli or categorical)
 - more flexible but require more data to avoid overfitting
- If the data is truly from a Naive Bayes model, then Logistic Regression and Naive Bayes converge to the exact same result in the limit (but NB is faster).

Today

Logistic regression

- Modeling
- Learning
- Learning (implementation)
- GDA vs. NB vs. Logistic regression
- Another perspective



Logistic regression: another perspective

- **Training data:** $D_{tr} = \{ (x_i \in \mathbb{R}^D, y_i \in \{+1, -1\}) \}_{i=1}^N$
- Learning: $\tilde{\mathbf{w}}^{\text{MLE}} = \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_{i=1}^N \log \left(1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{x}_i} \right)$
- Prediction: $\operatorname{sign}(\mathbf{w}^T \mathbf{x} + b)$

Logistic regression: another perspective

- Training data: $D_{tr} = \{ (x_i \in \mathbb{R}^D, y_i \in \{+1, -1\}) \}_{i=1}^N$
 - Learning: $\tilde{\mathbf{w}}^{\text{MLE}} = \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_{i=1}^N \log \left(1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{x}_i} \right) =$
$$\underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_{i=1}^N \ell(\tilde{\mathbf{w}}^T \tilde{x}_i, y_i) = \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_{i=1}^N \ell(h_{\tilde{\mathbf{w}}}(\tilde{x}_i), y_i) =$$
$$\underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_{i=1}^N \ell(\hat{y}_i, y_i)$$
 - Prediction: $\operatorname{sign}(\mathbf{w}^T \mathbf{x} + b)$
- “soft” prediction! Loss function

Empirical risk minimization recap

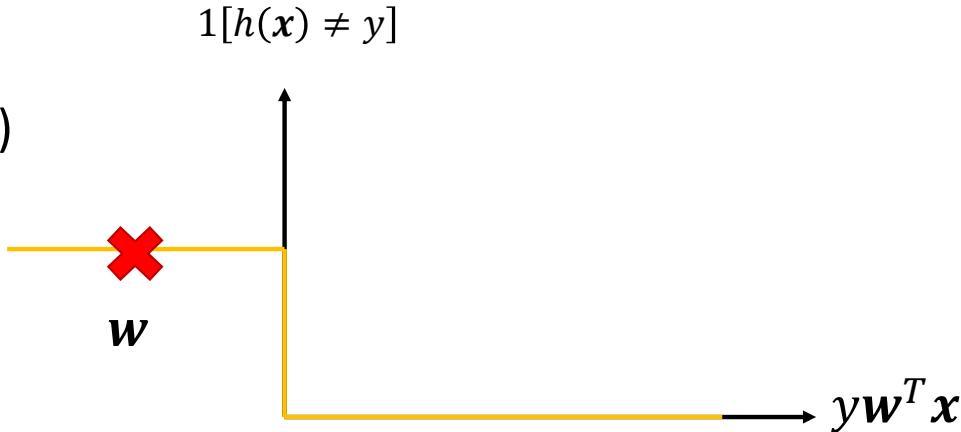
- In order to achieve a low generalization loss $E_{(x,y) \sim P}[\ell(h(x), y)]$
 - We minimize the training loss + certain regularization
 - Training loss: $\frac{1}{N} \sum_{n=1}^N \ell(h(x_n), y_n)$, where $D_{tr} = \{(x_n \in \mathbb{R}^D, y_n \in \{+1, -1\})\}_{i=1}^N$
 - Ideally, h & ℓ in the training should be the same as h & ℓ in the test (generalization) loss

Empirical risk minimization recap

- In order to achieve a low generalization loss $E_{(x,y) \sim P}[\ell(h(x), y)]$
 - We minimize the training loss + certain regularization
 - Training loss: $\frac{1}{N} \sum_{n=1}^N \ell(h(x_n), y_n)$, where $D_{tr} = \{(x_n \in \mathbb{R}^D, y_n \in \{+1, -1\})\}_{i=1}^N$
 - Ideally, h & ℓ in the training should be the same as h & ℓ in the test (generalization) loss
- For classification, we care
 - 0-1 loss: $\ell(y, \hat{y}) = 1[y \neq \hat{y}]$; $1[\text{True}] = 1, 1[\text{False}] = 0$ is called an indicator function
- For linear classifier, we use
 - $\hat{y} = h(x) = \text{sign}(w^T x)$; the bias term b is absorbed into w for brevity

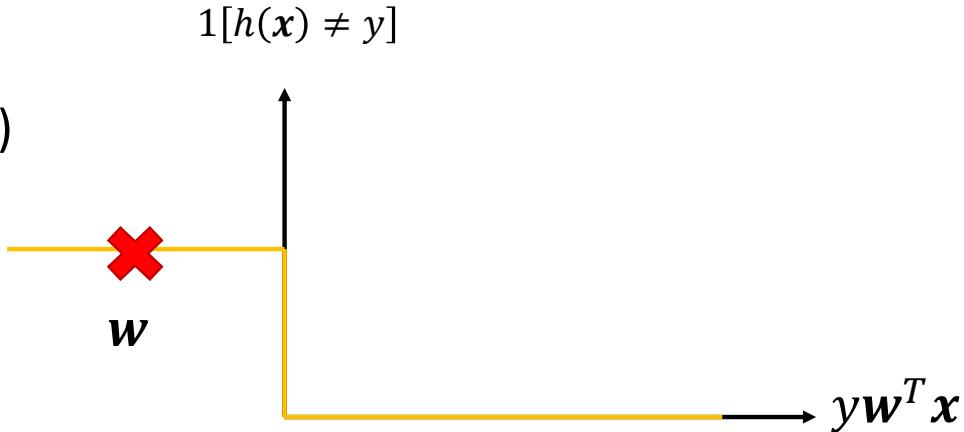
Problems with 0-1 losses

- $\ell(y, \hat{y}) = 1[y \neq \hat{y}]$
 - Not differentiable w.r.t. \hat{y} (i.e., the prediction)
 - Example: for $y_i \in \{+1, -1\}$, we want
 - \hat{y} (and thus $w^T x$) and y to be of the same sign



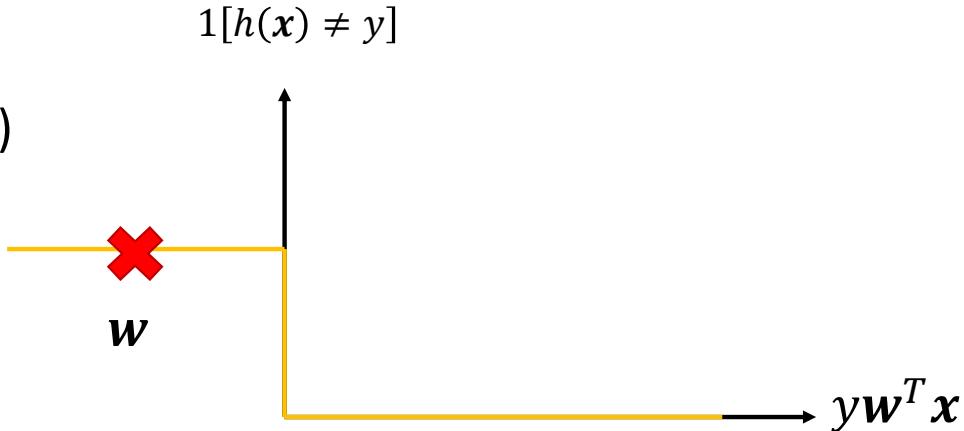
Problems with 0-1 losses

- $\ell(y, \hat{y}) = 1[y \neq \hat{y}]$
 - Not differentiable w.r.t. \hat{y} (i.e., the prediction)
 - Example: for $y_i \in \{+1, -1\}$, we want
 - \hat{y} (and thus $w^T x$) and y to be of the same sign
 - Hard to optimize w
 - No signals of where w should move



Problems with 0-1 losses

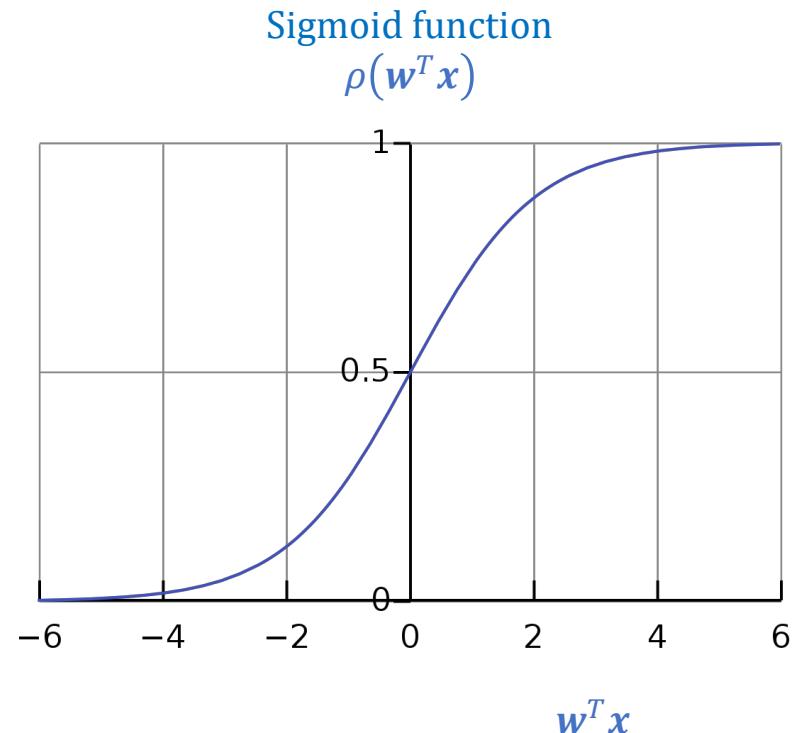
- $\ell(y, \hat{y}) = 1[y \neq \hat{y}]$
 - Not differentiable w.r.t. \hat{y} (i.e., the prediction)
 - Example: for $y_i \in \{+1, -1\}$, we want
 - \hat{y} (and thus $w^T x$) and y to be of the same sign
 - Hard to optimize w
 - No signals of where w should move



- We need a relaxed (surrogate) loss “**in training**” to encourage:
 - $y = 1, w^T x \uparrow$
 - $y = -1, w^T x \downarrow$

Logistic loss

- Let $p(+1|x; \mathbf{w}) = \rho(\mathbf{w}^T \mathbf{x})$
- Let $p(-1|x; \mathbf{w}) = 1 - \rho(\mathbf{w}^T \mathbf{x})$
- Want to maximize: $p(y_n|x_n; \mathbf{w})$
- Want to minimize: $-\log p(y_n|x_n; \mathbf{w})$
 $= -\log \rho(y_n \mathbf{w}^T \mathbf{x})$
- When
 - $y = 1, \mathbf{w}^T \mathbf{x} \uparrow$
 - $y = -1, \mathbf{w}^T \mathbf{x} \downarrow$

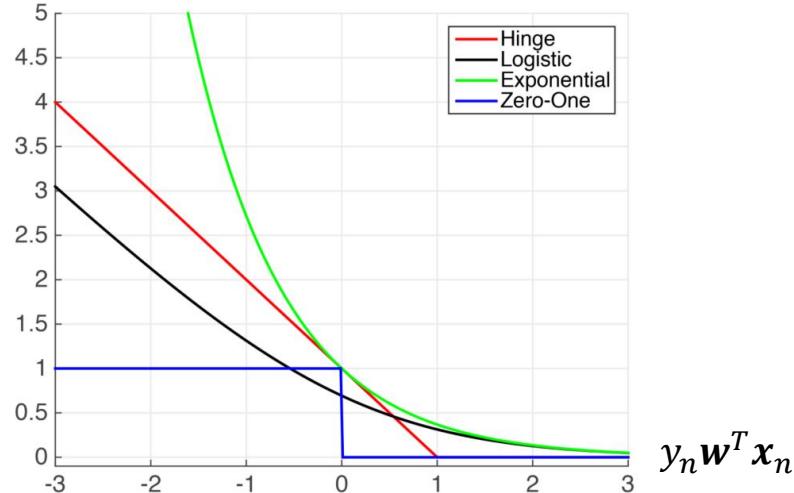


Questions!

- What properties must a surrogate loss have?
- Are there any other loss functions?
- Are there any other regularization functions, beyond $\lambda\|w\|_2^2$?

A first glance

- What properties must a surrogate loss have?
 - Upper bounding the 0-1 loss: so, small surrogate loss implies small 0-1 loss
- Are there any other loss functions?



- Are there any other regularization functions, beyond $\lambda \|\mathbf{w}\|_2^2$?
 - Yes: as long as it is “lower bounded” and characterizes model complexity

CSE 5523: Linear SVM



THE OHIO STATE UNIVERSITY

Logistic regression

- MLE: $\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(D_{tr}; \boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_n p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) p(\mathbf{x}_n)$

$$\begin{aligned} \log \prod_{i=1}^N p(y_i | \mathbf{x}_i; \boldsymbol{\theta}) &= \sum_{i=1}^N \log p(y_i | \mathbf{x}_i; \boldsymbol{\theta}) = \sum_{i=1}^N \log \rho(y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i) \\ &= \sum_{i=1}^N \log \frac{1}{1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i}} = \sum_{i=1}^N -\log(1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}_i}) \end{aligned}$$

- How to maximize?
 - Gradient descent!

GDA vs. NB vs. Logistic regression

- Training data: $D_{tr} = \{ (\mathbf{x}_i \in \mathbb{R}^D, y_i \in \{+1, -1\}) \}_{i=1}^N$
- Under certain assumptions, all of them will lead to a linear classifier $\text{sign}(\mathbf{w}^T \mathbf{x} + b)$
- However, the learned \mathbf{w} and b are different
 - GDA and NB need NO gradient descents; they have closed-form solutions
 - GDA and NB are suitable for streaming settings

GDA vs. NB vs. Logistic regression

- Training data: $D_{tr} = \{ (\mathbf{x}_i \in \mathbb{R}^D, y_i \in \{+1, -1\}) \}_{i=1}^N$
- Under certain assumptions, all of them will lead to a linear classifier $\text{sign}(\mathbf{w}^T \mathbf{x} + b)$
- However, the learned \mathbf{w} and b are different
 - GDA and NB need NO gradient descents; they have closed-form solutions
 - GDA and NB are suitable for streaming settings
- **Recap:**
 - A learning algorithm + A hypothesis class
 - GDA, NB, and logistic regression use the same hypothesis class but different algorithms

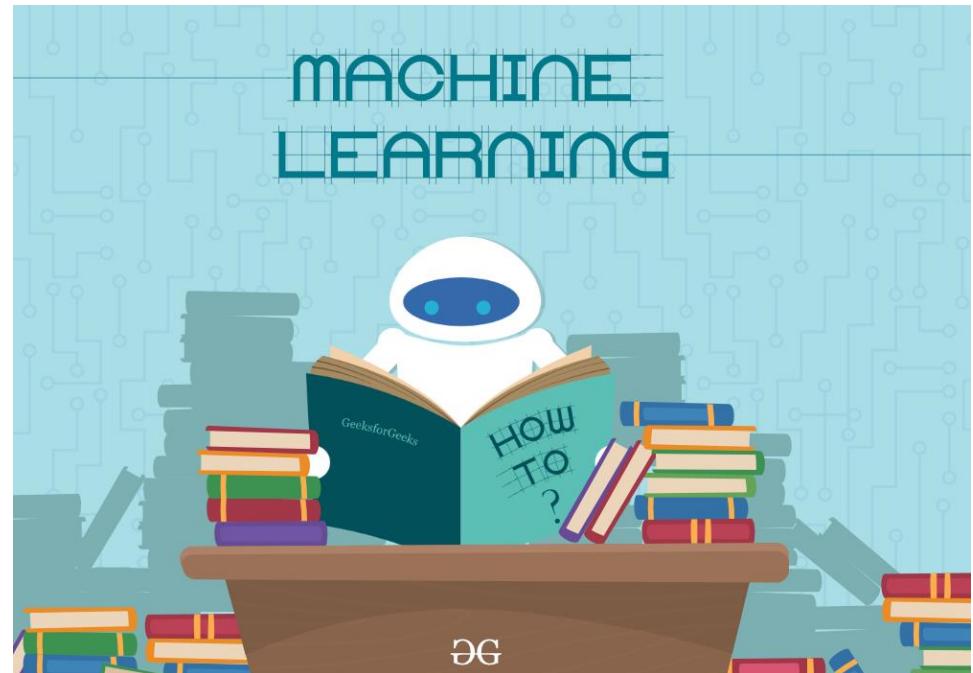
Generative vs. Discriminative learning

- Generative learning (MLE case)
 - $p(x, y; \theta) = p(x|y; \theta)p(y; \theta)$
 - Make **two modeling assumptions** (conditional data + Bernoulli or categorical)
 - If the assumption is correct, then fewer training data are needed
- Discriminative learning (MLE case)
 - $p(x, y; \theta) = p(y|x; \theta)p(x)$
 - Make **one modeling assumptions** (conditional Bernoulli or categorical)
 - more flexible but require more data to avoid overfitting
- If the data is truly from a Naive Bayes model, then Logistic Regression and Naive Bayes converge to the exact same result in the limit (but NB is faster).

Today

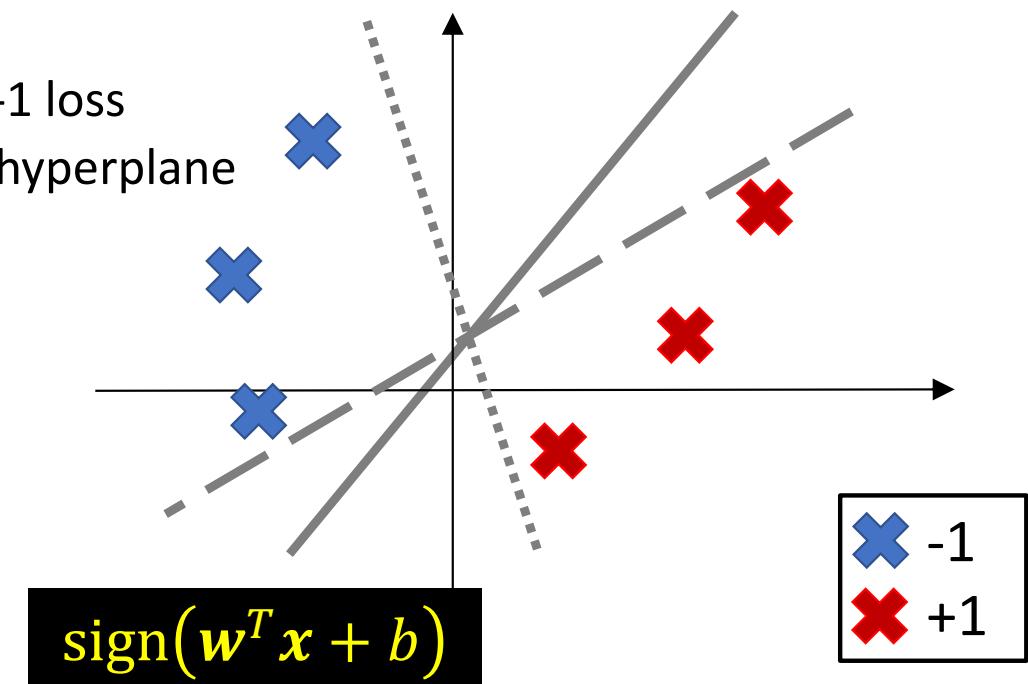
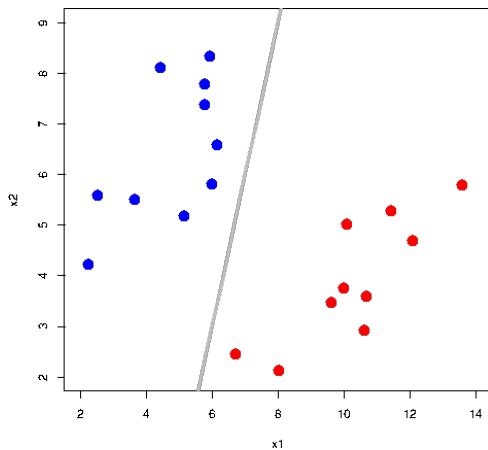
Support vector machines (SVM)

- Overview
- Margin
- Maximum margin classifier
- Support vectors
- Soft-margin SVM
- Summary



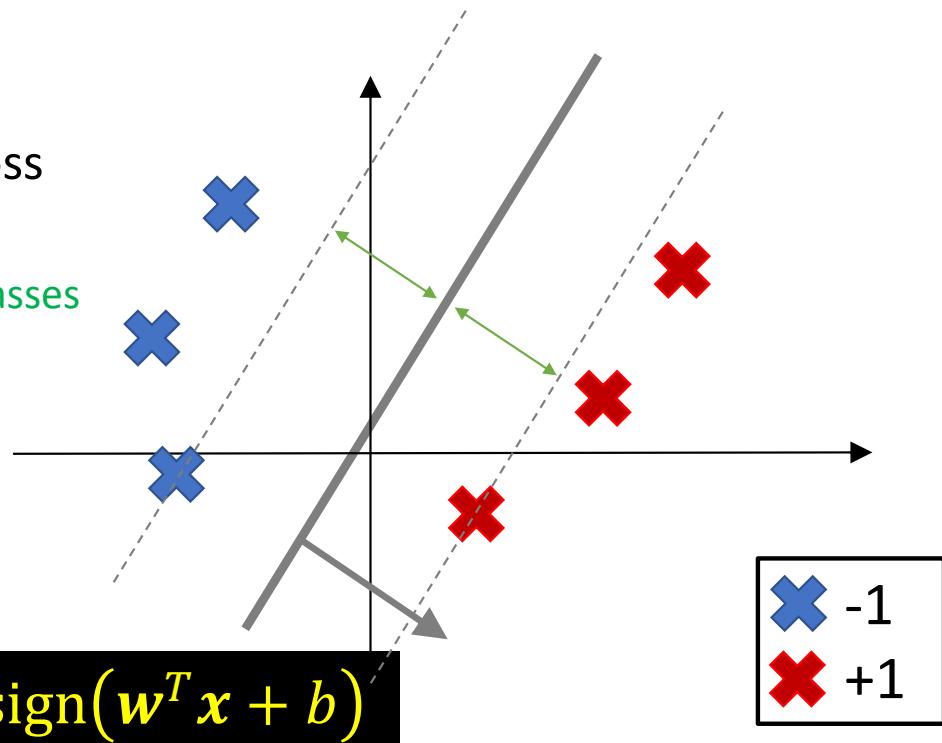
Support vector machines (SVM)

- A linear classifier viewed as an extension of the perceptron algorithm
- Given a linear separable D_{tr}
 - There are infinite many (w, b) for zero 0-1 loss
 - Perceptron finds an arbitrary separating hyperplane



Support vector machines (SVM)

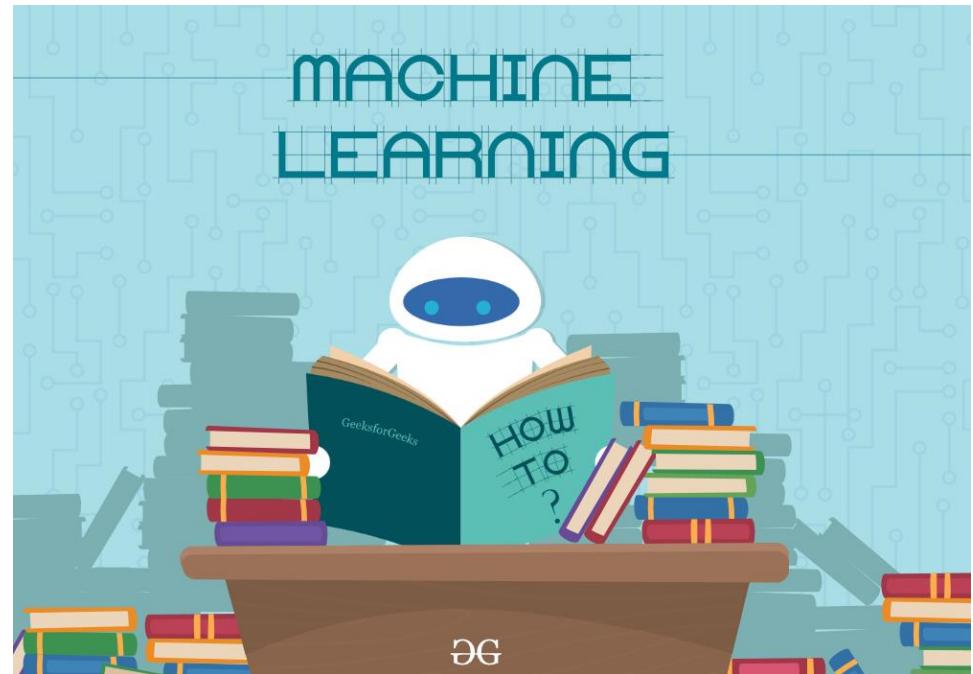
- A linear classifier viewed as an extension of the perceptron algorithm
- Given a linear separable D_{tr}
 - There are infinite many (w, b) for zero 0-1 loss
 - SVM picks the one that maximizes
 - the distance to the closest data points of both classes
 - SVM finds the maximum margin hyperplane



Today

Support vector machines (SVM)

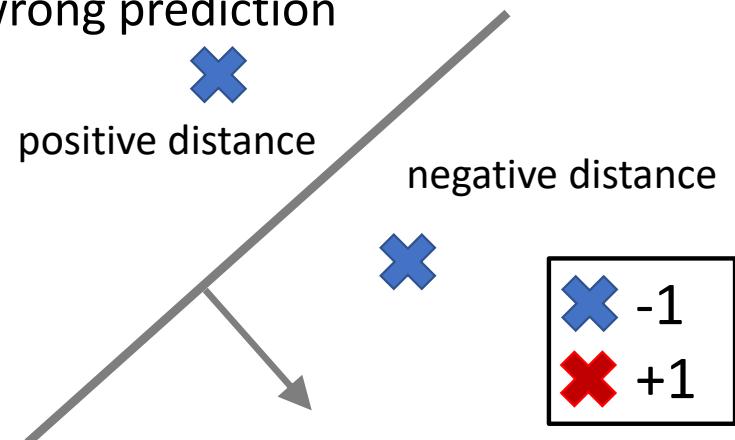
- Overview
- Margin
- Maximum margin classifier
- Support vectors
- Soft-margin SVM
- Summary



Margin

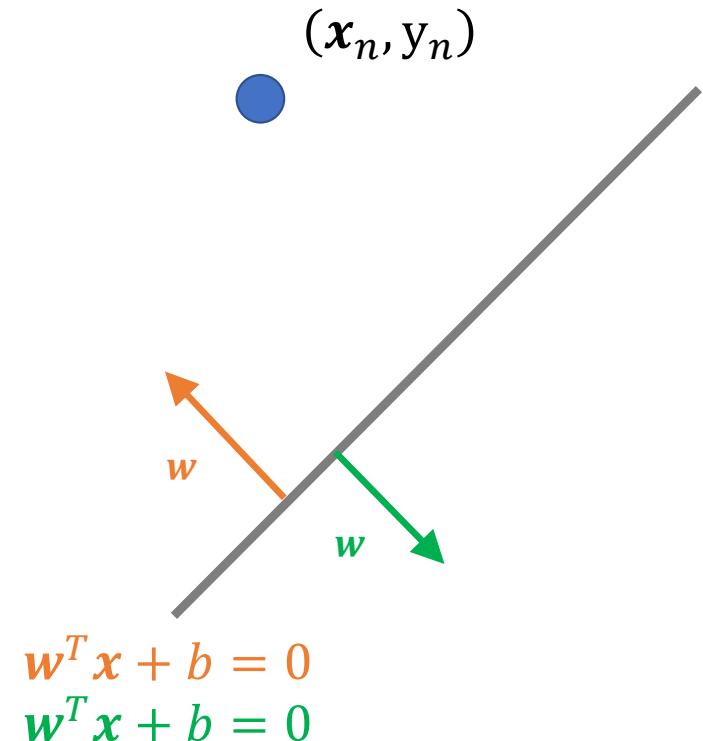
- Hypothesis class: hyperplane (a.k.a. linear boundary)
 - $\mathcal{H} = \{\mathbf{w}^T \mathbf{x} + b = 0\}$; here I explicitly write down b
- Given $D_{tr} = \{(x_n \in \mathbb{R}^D, y_n \in \{+1, -1\})\}_{i=1}^N$
 - Margin γ = “sign” distance from the hyperplane to the closest points of both classes
 - Sign: “positive” for correct prediction; “negative” for wrong prediction

$$\text{sign}(\mathbf{w}^T \mathbf{x} + b)$$



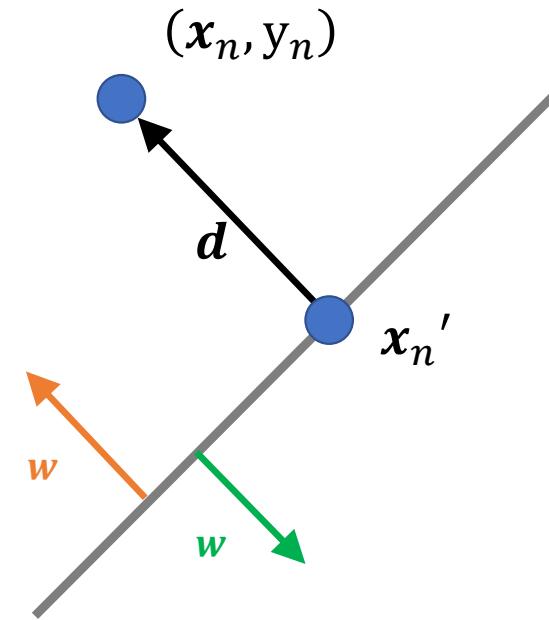
Margin

- Let's begin with one data point (x_n, y_n) , what is its distance to a hyperplane?



Margin

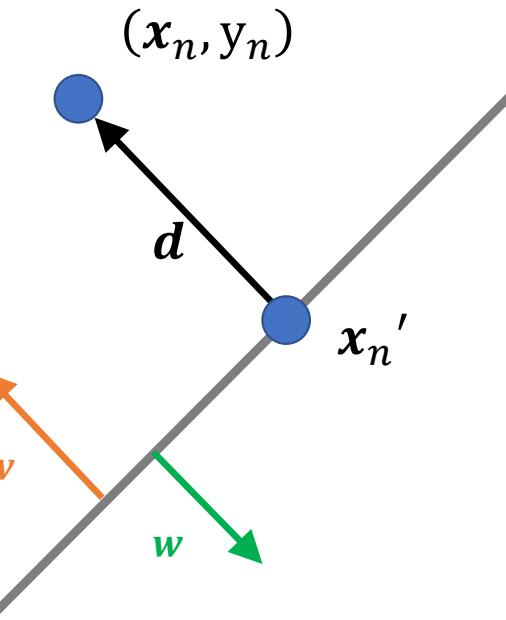
- Let's begin with one data point (x_n, y_n) , what is its distance to a hyperplane?
 - Let x_n' be x_n projection onto the hyperplane
 - Let d be the vector such that $x_n' + d = x_n$



$$w^T x + b = 0$$
$$w^T x + b = 0$$

Margin

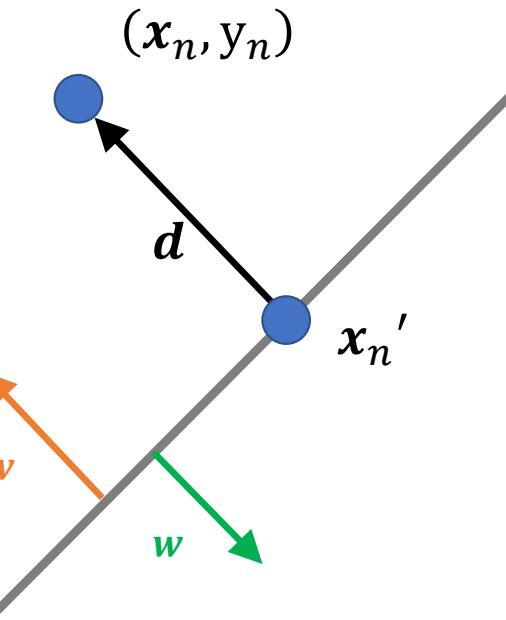
- Let's begin with one data point (x_n, y_n) , what is its distance to a hyperplane?
 - Let x_n' be x_n projection onto the hyperplane
 - Let d be the vector such that $x_n' + d = x_n$
- Some properties
 - $d = \alpha w$, where $\alpha \in \mathbb{R}$ (i.e., in parallel)
 - $w^T x_n' + b = 0$



$$w^T x + b = 0$$
$$w^T x + b = 0$$

Margin

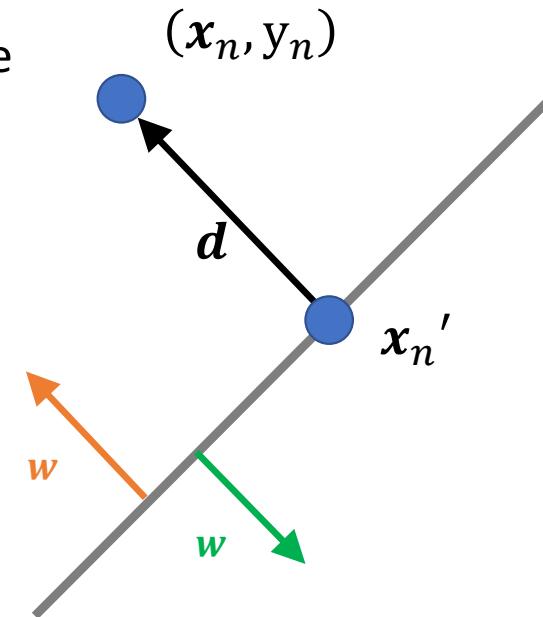
- Let's begin with one data point (x_n, y_n) , what is its distance to a hyperplane?
 - Let x_n' be x_n projection onto the hyperplane
 - Let d be the vector such that $x_n' + d = x_n$
- Some properties
 - $d = \alpha w$, where $\alpha \in \mathbb{R}$ (i.e., in parallel)
 - $w^T x_n' + b = 0$
- $w^T x_n' + b = w^T(x_n - d) + b = w^T(x_n - \alpha w) + b = 0$
 - $\alpha = \frac{w^T x_n + b}{w^T w}$
 - $\|d\|_2 = \sqrt{d^T d} = \sqrt{\alpha^2 w^T w} = \frac{|w^T x_n + b|}{\|w\|_2}$



$$w^T x + b = 0$$
$$w^T x + b = 0$$

Margin

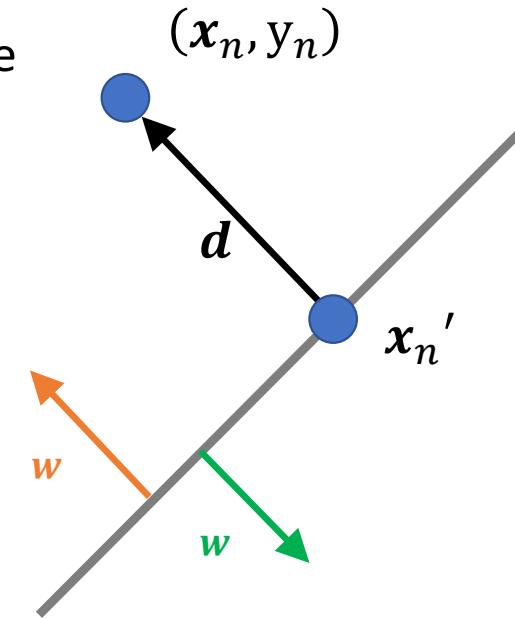
- “Sign” distance for d : $\frac{y_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|_2}$
 - Non-negative if (\mathbf{x}_n, y_n) is on the correct side of the hyperplane
 - For $\mathbf{w}^T \mathbf{x} + b = 0$, non-negative if $y_n = 1$
 - For $\mathbf{w}^T \mathbf{x} + b = 0$, non-negative if $y_n = -1$



$$\begin{aligned}\mathbf{w}^T \mathbf{x} + b &= 0 \\ \mathbf{w}^T \mathbf{x} + b &= 0\end{aligned}$$

Margin

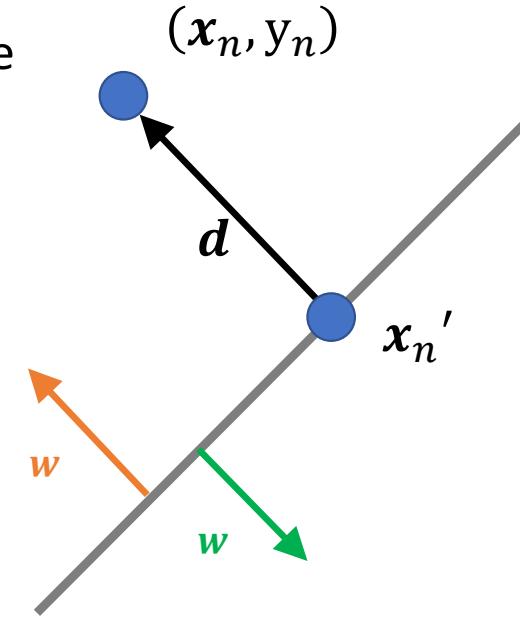
- “Sign” distance for d : $\frac{y_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|_2}$
 - Non-negative if (\mathbf{x}_n, y_n) is on the correct side of the hyperplane
 - For $\mathbf{w}^T \mathbf{x} + b = 0$, non-negative if $y_n = 1$
 - For $\mathbf{w}^T \mathbf{x} + b = 0$, non-negative if $y_n = -1$
- Properties:
 - The hyperplane and distance are positive-scale-invariant
 - $\frac{y_n(\beta \mathbf{w}^T \mathbf{x}_n + \beta b)}{\|\beta \mathbf{w}\|_2} = \frac{y_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|_2}$ for $\beta > 0$



$$\begin{aligned}\mathbf{w}^T \mathbf{x} + b &= 0 \\ \mathbf{w}^T \mathbf{x} + b &= 0\end{aligned}$$

Margin

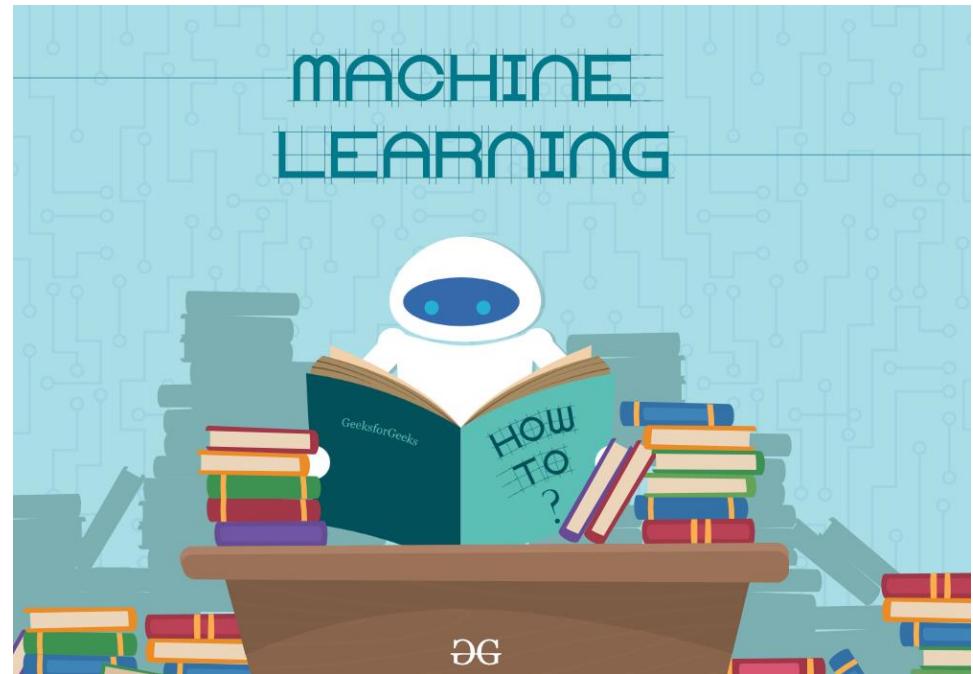
- “Sign” distance for d : $\frac{y_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|_2}$
 - Non-negative if (\mathbf{x}_n, y_n) is on the correct side of the hyperplane
 - For $\mathbf{w}^T \mathbf{x} + b = 0$, non-negative if $y_n = 1$
 - For $\mathbf{w}^T \mathbf{x} + b = 0$, non-negative if $y_n = -1$
- Properties:
 - The hyperplane and distance are positive-scale-invariant
 - $\frac{y_n(\beta \mathbf{w}^T \mathbf{x}_n + \beta b)}{\|\beta \mathbf{w}\|_2} = \frac{y_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|_2}$ for $\beta > 0$
- Margin for D_{tr} : $\gamma = \min_{(\mathbf{x}, y) \in D_{tr}} \frac{y(\mathbf{w}^T \mathbf{x} + b)}{\|\mathbf{w}\|_2}$



Today

Support vector machines (SVM)

- Overview
- Margin
- Maximum margin classifier
- Support vectors
- Soft-margin SVM
- Summary



Maximum margin classifier

- Training data: $D_{tr} = \{ (x_n \in \mathbb{R}^D, y_n \in \{+1, -1\}) \}_{i=1}^N$ (linear separable)
- $(w, b)^{\text{SVM}} = \operatorname{argmax}_{(w, b)} \gamma(w, b) = \operatorname{argmax}_{(w, b)} \min_{(x_n, y_n)} \frac{y_n(w^T x_n + b)}{\|w\|_2} = \operatorname{argmax}_{(w, b)} \frac{1}{\|w\|_2} \min_{(x_n, y_n)} y_n(w^T x_n + b)$

Maximum margin classifier

- Training data: $D_{tr} = \{ (\mathbf{x}_n \in \mathbb{R}^D, y_n \in \{+1, -1\}) \}_{i=1}^N$ (linear separable)

- $(\mathbf{w}, b)^{\text{SVM}} = \operatorname{argmax}_{(\mathbf{w}, b)} \gamma(\mathbf{w}, b) = \operatorname{argmax}_{(\mathbf{w}, b)} \min_{(\mathbf{x}_n, y_n)} \frac{y_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|_2} = \operatorname{argmax}_{(\mathbf{w}, b)} \frac{1}{\|\mathbf{w}\|_2} \min_{(\mathbf{x}_n, y_n)} y_n(\mathbf{w}^T \mathbf{x}_n + b)$



independent of (\mathbf{x}_n, y_n)

Maximum margin classifier

- Training data: $D_{tr} = \{(\mathbf{x}_n \in \mathbb{R}^D, y_n \in \{+1, -1\})\}_{i=1}^N$ (linear separable)
- $(\mathbf{w}, b)^{\text{SVM}} = \operatorname{argmax}_{(\mathbf{w}, b)} \gamma(\mathbf{w}, b) = \operatorname{argmax}_{(\mathbf{w}, b)} \min_{(\mathbf{x}_n, y_n)} \frac{y_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|_2} = \operatorname{argmax}_{(\mathbf{w}, b)} \frac{1}{\|\mathbf{w}\|_2} \min_{(\mathbf{x}_n, y_n)} y_n(\mathbf{w}^T \mathbf{x}_n + b)$
- Due to **positive-scale-invariant**, for (\mathbf{w}, b) that perfectly separates data, we can
 - constrain $\min_{(\mathbf{x}_n, y_n)} y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$
 - $(\mathbf{w}, b)^{\text{SVM}} = \operatorname{argmax}_{(\mathbf{w}, b)} \frac{1}{\|\mathbf{w}\|_2}$, s.t. $\min_{(\mathbf{x}_n, y_n)} y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$

Maximum margin classifier

- $(\mathbf{w}, b)^{\text{SVM}} = \underset{(\mathbf{w}, b)}{\operatorname{argmax}} \frac{1}{\|\mathbf{w}\|_2} = \underset{(\mathbf{w}, b)}{\operatorname{argmin}} \|\mathbf{w}\|_2 = \underset{(\mathbf{w}, b)}{\operatorname{argmin}} \|\mathbf{w}\|_2^2 = \underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w}$
- Optimization problem:
 - $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w},$ s.t. $\min_{(\mathbf{x}_n, y_n)} y_n (\mathbf{w}^T \mathbf{x}_n + b) = 1$

Maximum margin classifier

- $(\mathbf{w}, b)^{\text{SVM}} = \operatorname{argmax}_{(\mathbf{w}, b)} \frac{1}{\|\mathbf{w}\|_2} = \operatorname{argmin}_{(\mathbf{w}, b)} \|\mathbf{w}\|_2 = \operatorname{argmin}_{(\mathbf{w}, b)} \|\mathbf{w}\|_2^2 = \operatorname{argmin}_{(\mathbf{w}, b)} \frac{1}{2} \mathbf{w}^T \mathbf{w}$
- Optimization problem:
 - $\operatorname{argmin}_{(\mathbf{w}, b)} \frac{1}{2} \mathbf{w}^T \mathbf{w}, \quad \text{s.t. } \min_{(\mathbf{x}_n, y_n)} y_n (\mathbf{w}^T \mathbf{x}_n + b) = 1$
- “Equivalent” optimization problem with the same solution
 - $\operatorname{argmin}_{(\mathbf{w}, b)} \frac{1}{2} \mathbf{w}^T \mathbf{w}, \quad \text{s.t. } \forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$
 - A quadratic optimization problem (quadratic objective; linear constraints) that can be solved efficiently via a QCQP (Quadratically Constrained Quadratic Program) solver.
 - To prove the equivalence in your next homework!

Maximum margin classifier

- Optimization problem for linear “hard-margin” SVM (primal form)

$$\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w}, \quad \text{s.t. } \forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$$

Maximum margin classifier

- Optimization problem for linear “hard-margin” SVM (primal form)

- $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w}, \quad \text{s.t. } \forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$

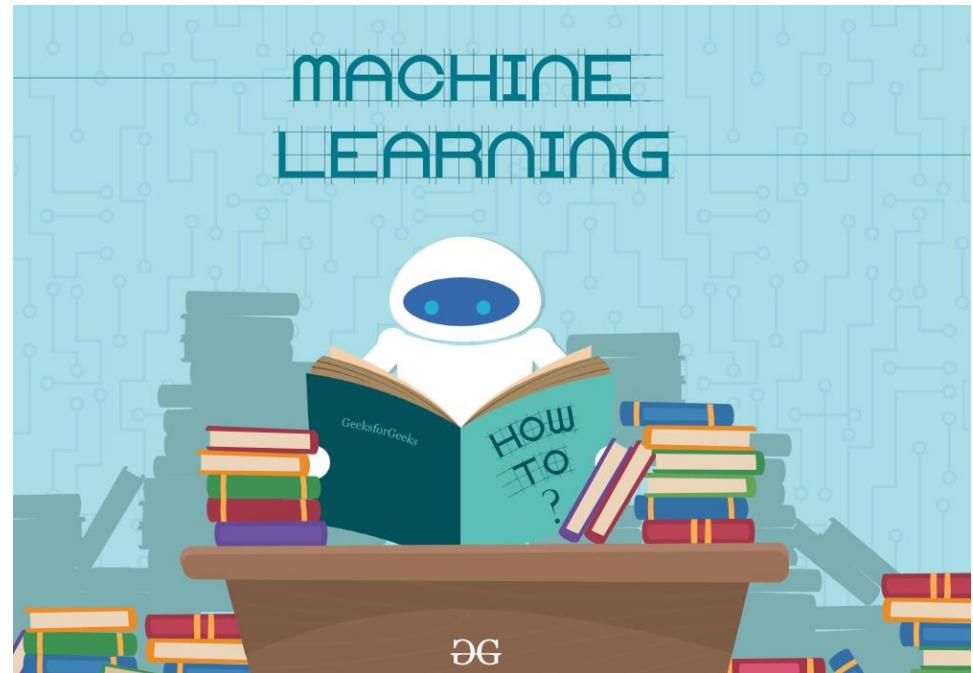
- Properties:

- Have a “unique” solution whenever a separating hyperplane exists
- $y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 0$ is already correct, but we ask more: $y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$
 - Perceptron only asks for the former $y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 0$
- Interpretation: find the simplest hyperplane (smallest $\|\mathbf{w}\|_2^2$) such that all data points lie at least one unit away from it
- Margin: $\frac{1}{\|\mathbf{w}^{\text{SVM}}\|_2}$

Today

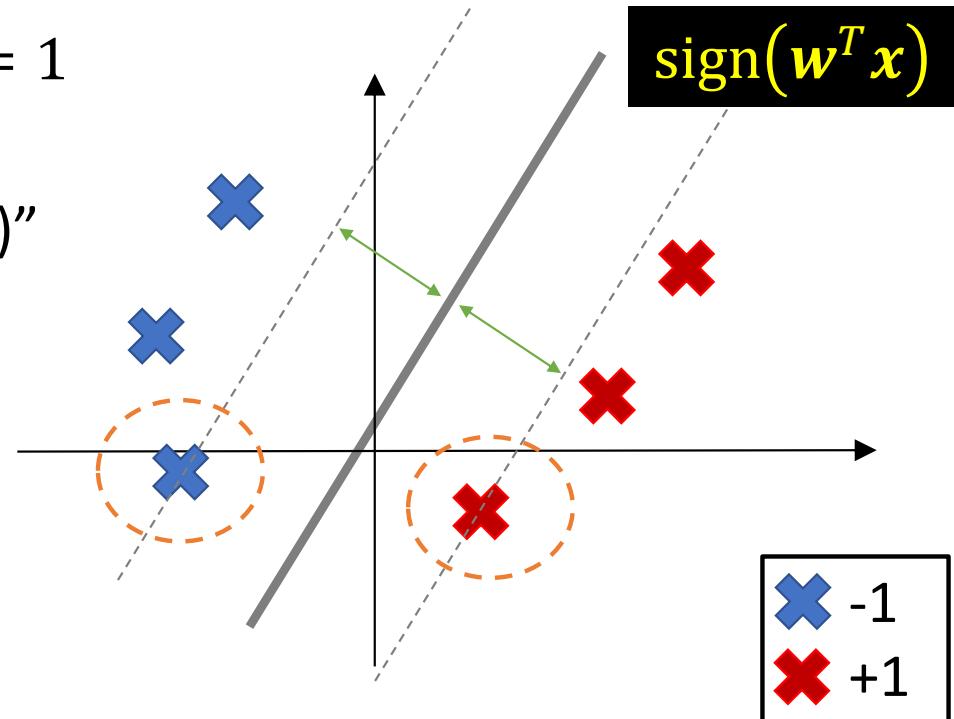
Support vector machines (SVM)

- Overview
- Margin
- Maximum margin classifier
- Support vectors
- Soft-margin SVM
- Summary



Support vectors

- Some (x_n, y_n) will have $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$
 - Homework!
- We refer to them as “support vectors (SV)”



Support vectors

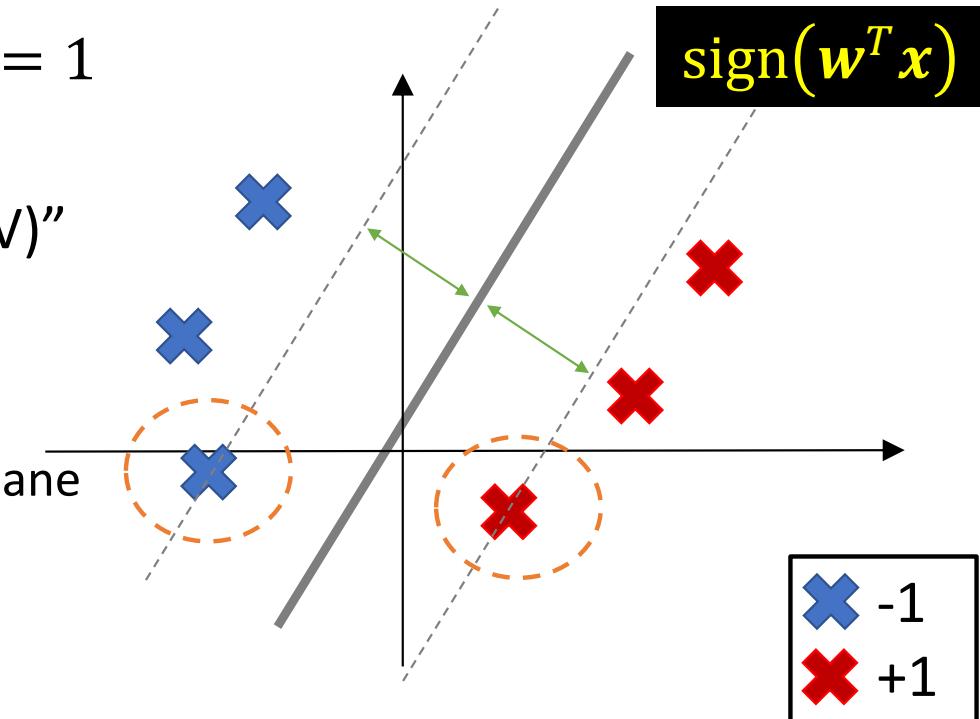
- Some (x_n, y_n) will have $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$
 - Homework!

- We refer to them as “support vectors (SV)”

- Properties:

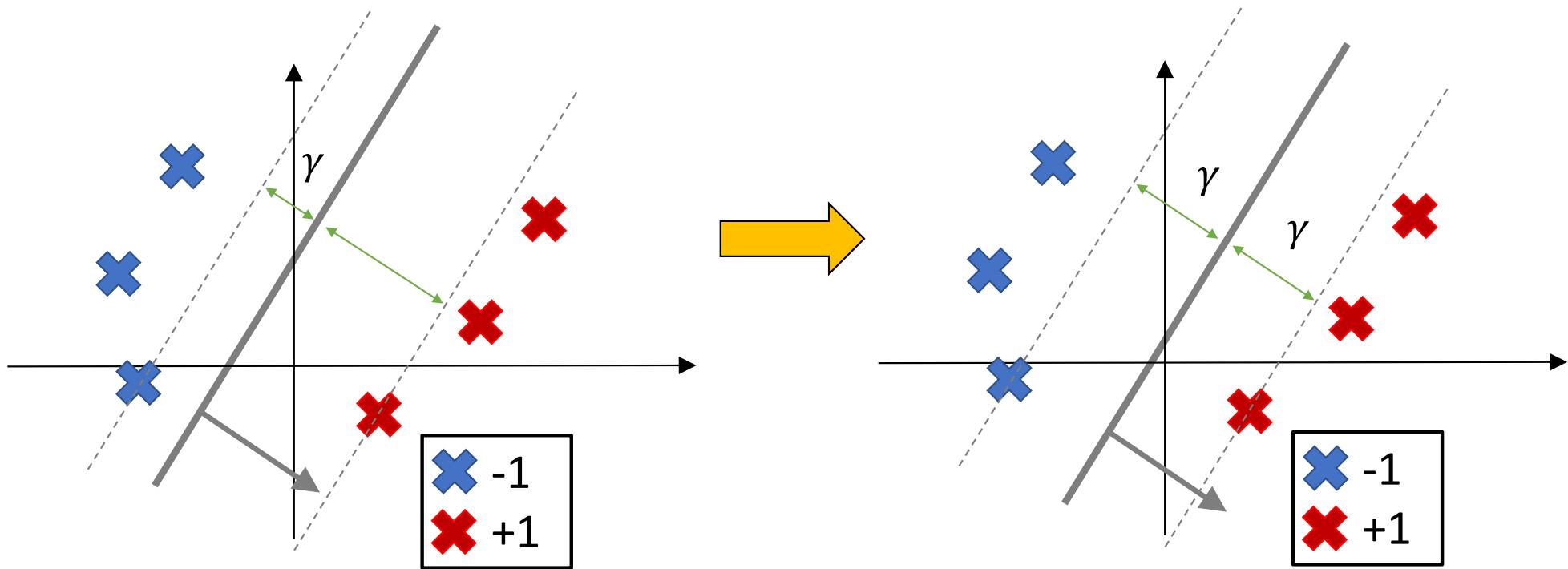
- SVs define the margin and thus the hyperplane
 - Removing them, $(\mathbf{w}, b)^{\text{SVM}}$ might change
 - **Removing other data points won't!**

- That is, you can obtain the same $(\mathbf{w}, b)^{\text{SVM}}$ from a (much) smaller D_{tr}



Properties

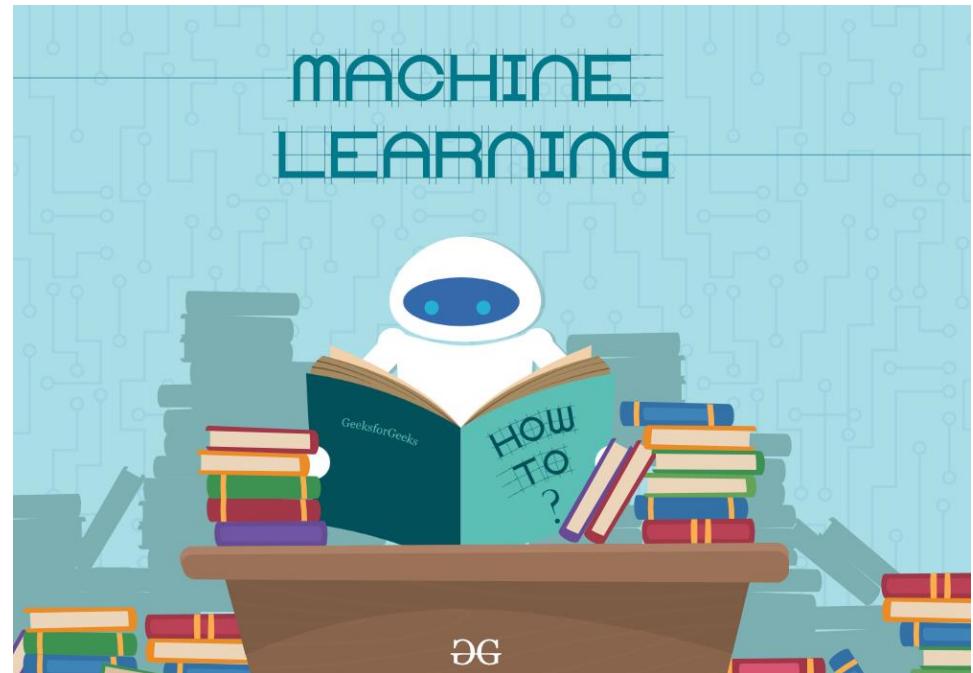
- A maximum-margin hyperplane lies in the middle of two classes



Today

Support vector machines (SVM)

- Overview
- Margin
- Maximum margin classifier
- Support vectors
- Soft-margin SVM
- Summary



Soft-margin SVM

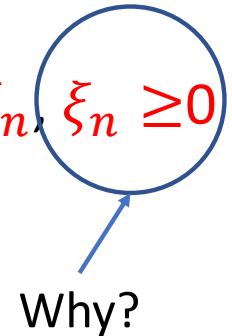
- $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w}, \quad \text{s.t. } \forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$

Soft-margin SVM

- $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w}$, s.t. $\forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$
 - No solution if the data is not linear separable (common cases)
- Allow “slight” violation with “slack” variables
- $\underset{(\mathbf{w}, b, \xi_n)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n \xi_n$, s.t. $\forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \xi_n \geq 0$
 - C is a hyper-parameter
 - How to choose it?

Soft-margin SVM

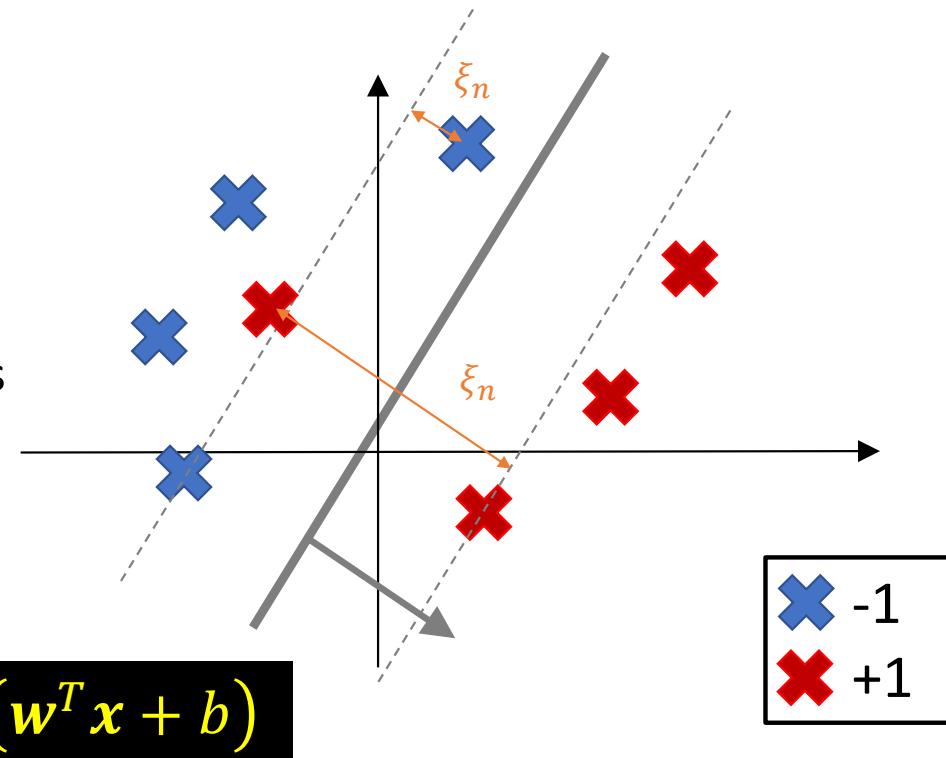
- $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w}$, s.t. $\forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$
 - No solution if the data is not linear separable (common cases)
- Allow “slight” violation with “slack” variables
- $\underset{(\mathbf{w}, b, \xi_n)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n \xi_n$, s.t. $\forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \xi_n \geq 0$
 - C is a hyper-parameter
 - How to choose it? **A: validation**



Why?

Slack variables

- Allow some data points to be closer to the hyperplane or classified wrongly
 - But with penalties in the objective
 - $\frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n \xi_n$
- Extreme cases:
 - If C is large, more like the hard-margin SVM
 - If C is small, favor simple models with errors



Hinge loss

- Original soft-margin SVM

- $\underset{(\mathbf{w}, b, \xi_n)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n \xi_n,$ s.t. $\forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \xi_n \geq 0$
- Convex quadratic programming: quadratic objective + linear constraints

Hinge loss

- Original soft-margin SVM

- $\underset{(\mathbf{w}, b, \xi_n)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n \xi_n$,
s.t. $\forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \xi_n \geq 0$
- Convex quadratic programming: quadratic objective + linear constraints

- Unconstrained formulation

- $y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n \Rightarrow \xi_n \geq 1 - y_n (\mathbf{w}^T \mathbf{x}_n + b)$
- $\xi_n = \max\{1 - y_n (\mathbf{w}^T \mathbf{x}_n + b), 0\} = [1 - y_n (\mathbf{w}^T \mathbf{x}_n + b)]_+$
- \geq disappears because we always want to minimize ξ_n

Hinge loss

- Original soft-margin SVM

- $\underset{(\mathbf{w}, b, \xi_n)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n \xi_n$,
s.t. $\forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \xi_n \geq 0$
- Convex quadratic programming: quadratic objective + linear constraints

- Unconstrained formulation

- $y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n \Rightarrow \xi_n \geq 1 - y_n (\mathbf{w}^T \mathbf{x}_n + b)$
- $\xi_n = \max\{1 - y_n (\mathbf{w}^T \mathbf{x}_n + b), 0\} = [1 - y_n (\mathbf{w}^T \mathbf{x}_n + b)]_+$
- \geq disappears because we always want to minimize ξ_n

- $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \left[\frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n [1 - y_n (\mathbf{w}^T \mathbf{x}_n + b)]_+ \right]$

training “hinge” loss
regularization

Soft-margin SVM (primal form)

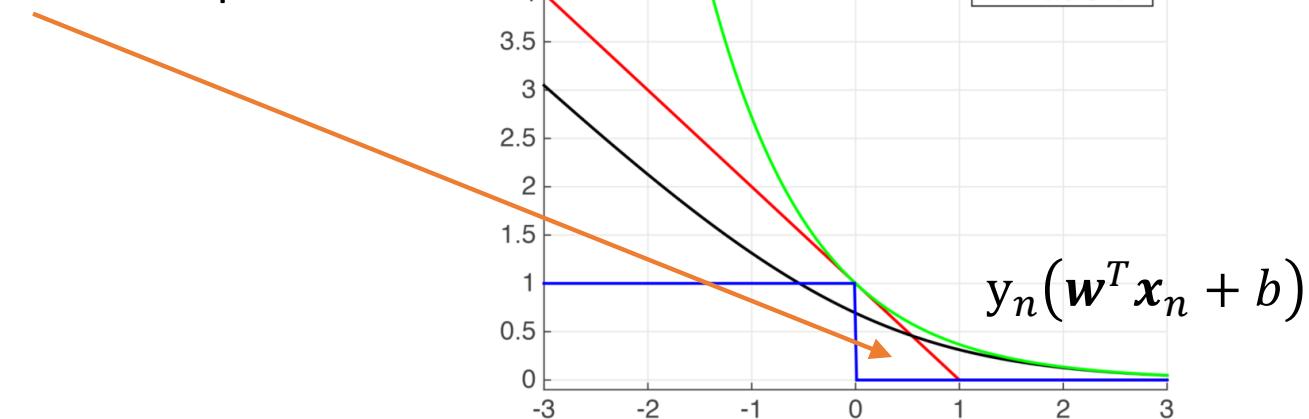
- Unconstrained problem: $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n [1 - y_n (\mathbf{w}^T \mathbf{x}_n + b)]_+$
 - How to solve?

Soft-margin SVM (primal form)

- Unconstrained problem: $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n [1 - y_n (\mathbf{w}^T \mathbf{x}_n + b)]_+$
 - How to solve? (stochastic) gradient descent
 - $[1 - y_n (\mathbf{w}^T \mathbf{x}_n + b)]_+$ is differentiable except at $1 = y_n (\mathbf{w}^T \mathbf{x}_n + b)$
 - Solve like logistic regression, but now hinge loss

Soft-margin SVM (primal form)

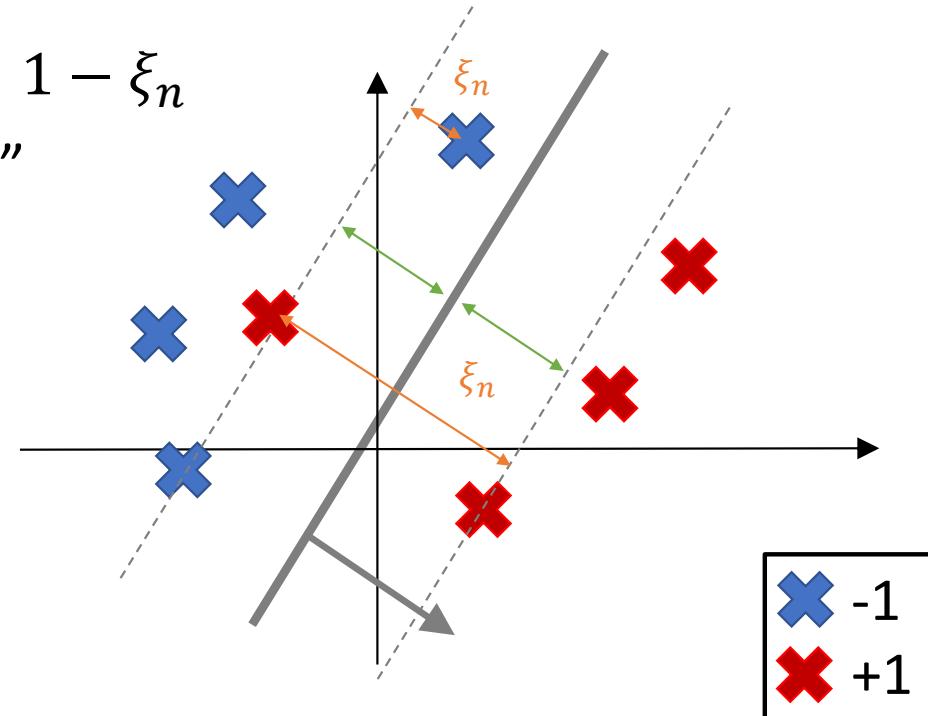
- Unconstrained problem: $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n [1 - y_n(\mathbf{w}^T \mathbf{x}_n + b)]_+$
 - How to solve? (stochastic) gradient descent
 - $[1 - y_n(\mathbf{w}^T \mathbf{x}_n + b)]_+$ is differentiable except at $1 = y_n(\mathbf{w}^T \mathbf{x}_n + b)$
 - Solve like logistic regression, but now hinge loss
 - Penalize even for unconfident correct predictions



Support vectors

$$\text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

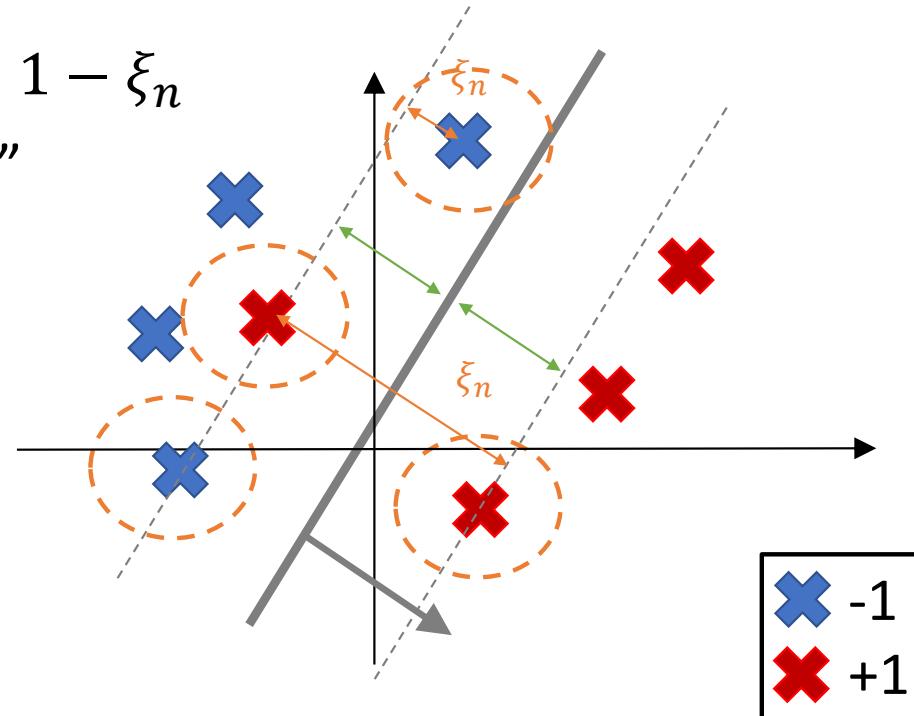
- Some (\mathbf{x}_n, y_n) will have $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1 - \xi_n$
- We refer to them as “support vectors (SV)”



Support vectors

$$\text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

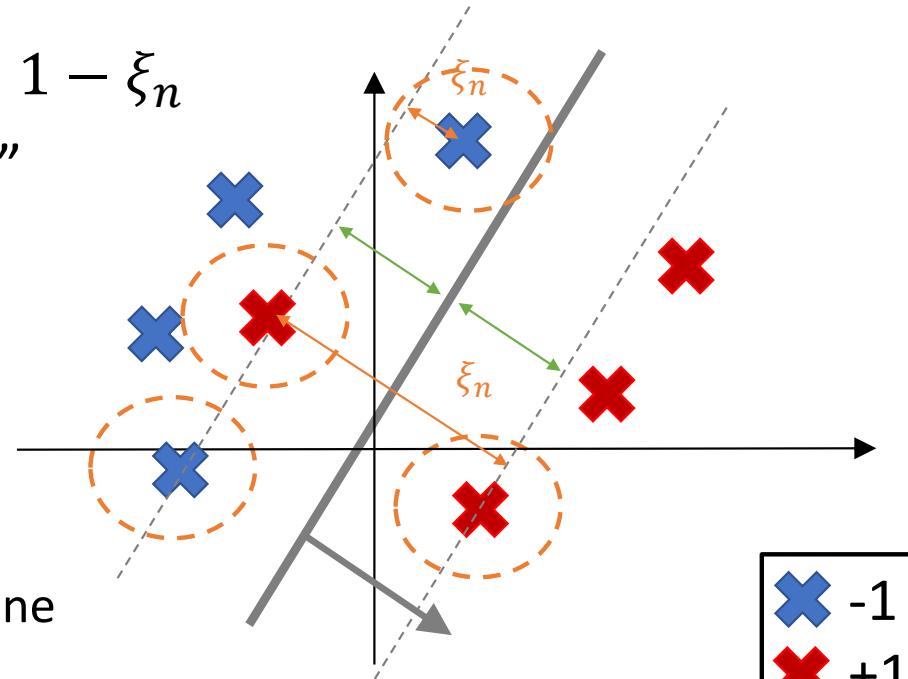
- Some (\mathbf{x}_n, y_n) will have $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1 - \xi_n$
- We refer to them as “support vectors (SV)”



Support vectors

$$\text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

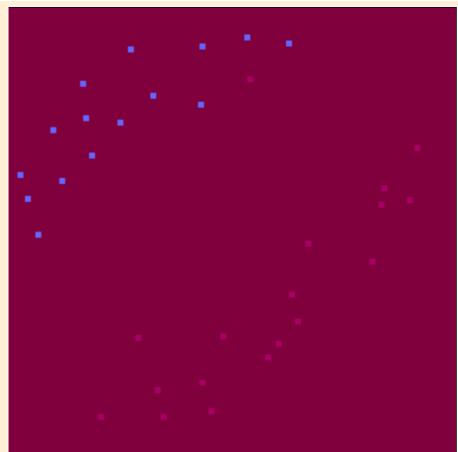
- Some (\mathbf{x}_n, y_n) will have $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1 - \xi_n$
- We refer to them as “support vectors (SV)”
- Properties: same as before
- Cases (for SVs):
 - $\xi_n = 0$: points are $\frac{1}{\|\mathbf{w}^{\text{SVM}}\|_2}$ from the hyperplane
 - $\xi_n < 1$: correctly predicted
 - $\xi_n > 1$: wrongly predicted



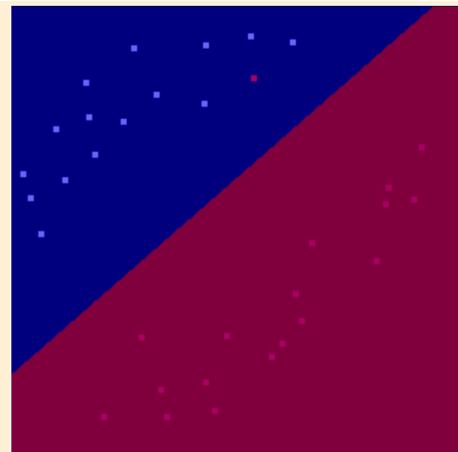
| | |
|--|----|
| | -1 |
| | +1 |

Demo

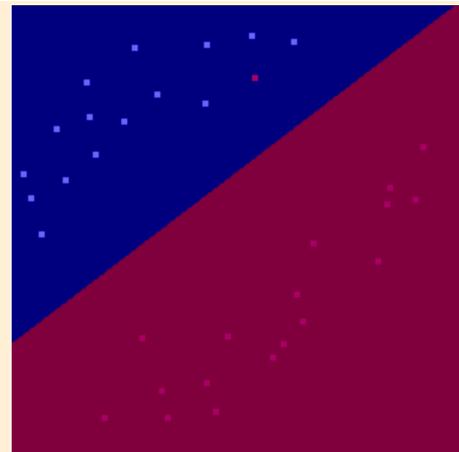
- LibSVM: <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>



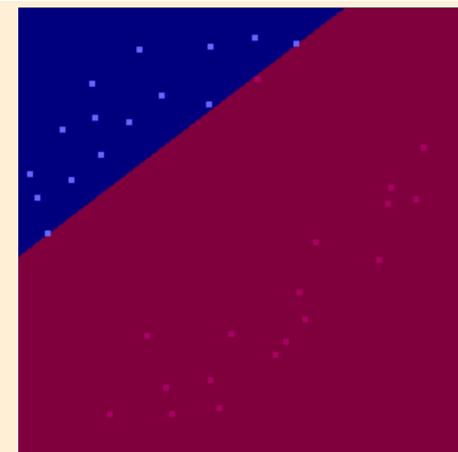
Change Run Clear -t 0 -c 0.01



Change Run Clear -t 0 -c 1



Change Run Clear -t 0 -c 100



Change Run Clear -t 0 -c 10000

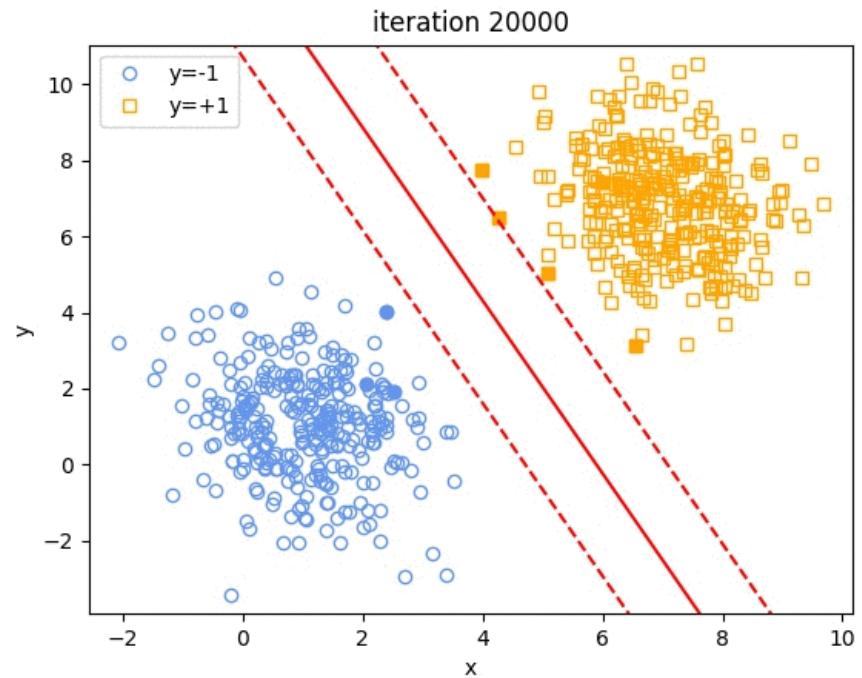
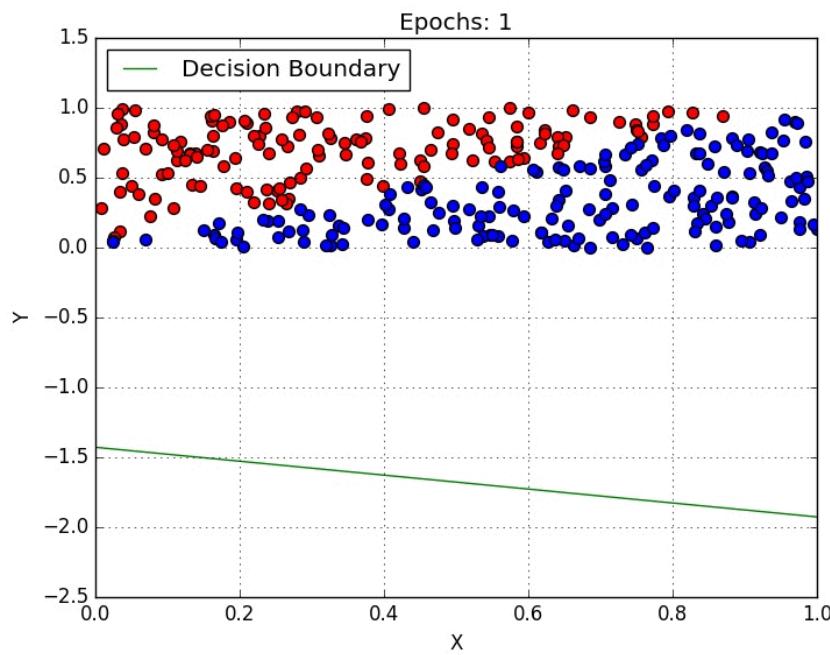
$C = 0.01$

$C = 1$

$C = 100$

$C = 10000$

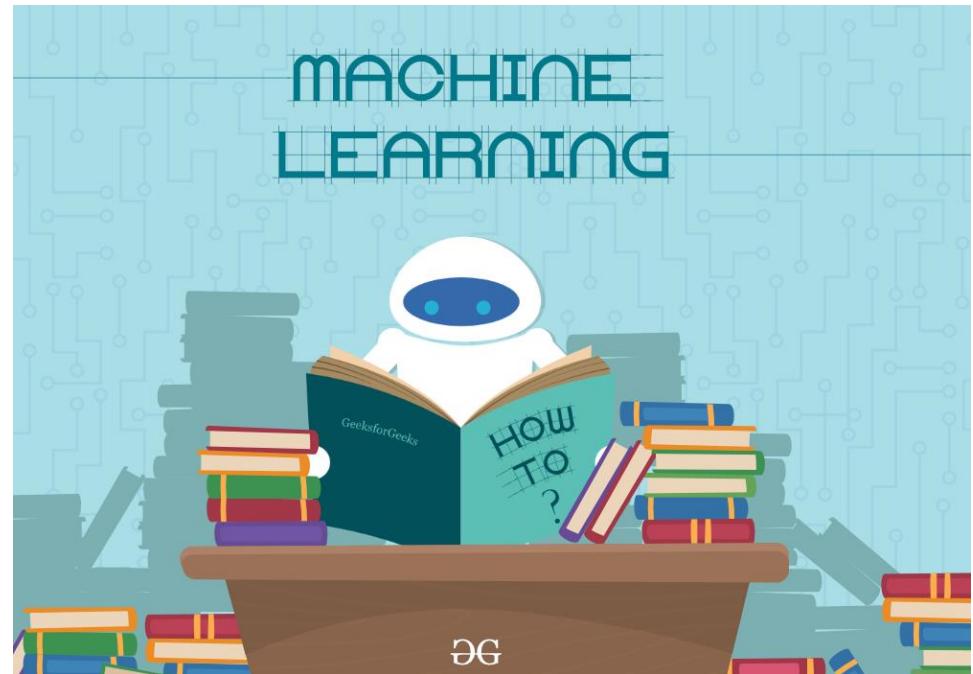
Iterative logistic regression and linear SVM



Today

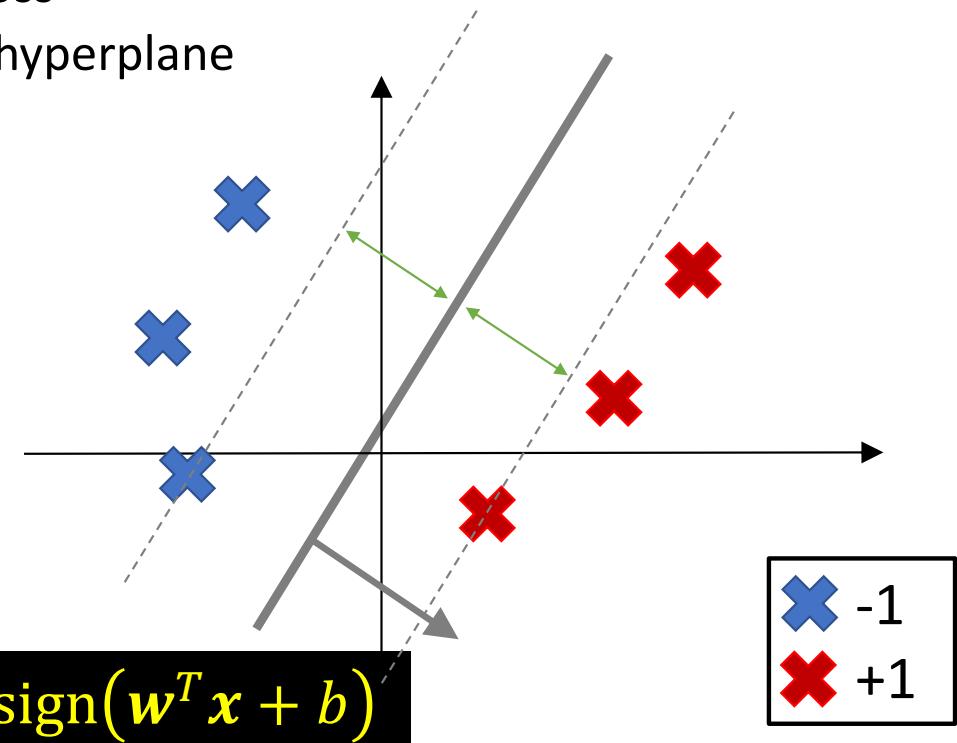
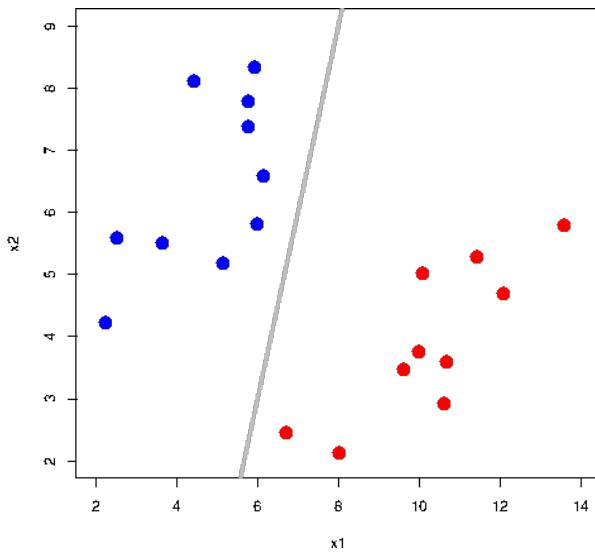
Support vector machines (SVM)

- Overview
- Margin
- Maximum margin classifier
- Support vectors
- Soft-margin SVM
- Summary



Summary: Support vector machines (SVM)

- Given a linear separable D_{tr}
 - There are infinite many (w, b) for zero 0-1 loss
 - SVM finds the maximum margin separating hyperplane



Summary: Hard-margin vs. Soft-margin SVM

- Linear “hard-margin” SVM

- $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w}$, s.t. $y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$

- No solution if the data is not linear separable (common cases)

- Linear “soft-margin” SVM

- $\underset{(\mathbf{w}, b, \xi_n)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n \xi_n$, s.t. $\forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \xi_n \geq 0$

- Allow “slight” violation with “slack” variables ξ_n

- C is a hyper-parameter

Summary: Hinge loss

- Original soft-margin SVM

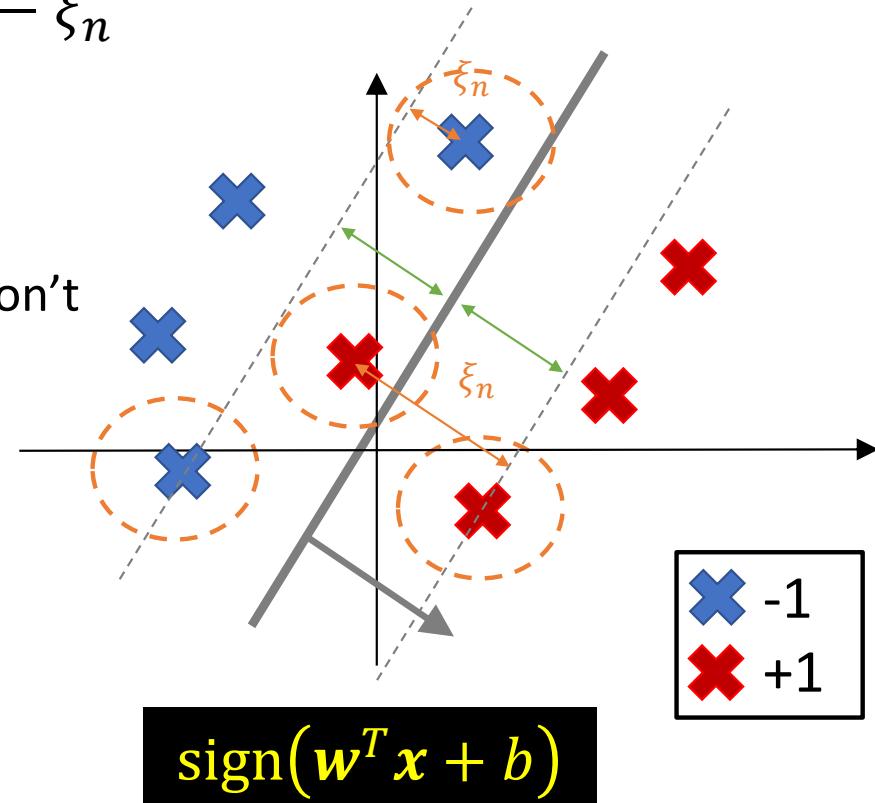
- $\underset{(\mathbf{w}, b, \xi_n)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n \xi_n,$ s.t. $\forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \xi_n \geq 0$
- Convex quadratic programming: quadratic objective + linear constraints
- Can be solved even faster by taking the special structures of the problem into account

- Unconstrained formulation

- $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n [1 - y_n (\mathbf{w}^T \mathbf{x}_n + b)]_+$ training “hinge” loss regularization

Summary: Support vectors

- Some (x_n, y_n) will have $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1 - \xi_n$
 - We refer to them as “support vectors (SV)”
 - SVs define the decision boundary
 - Removing them change the decision boundary
 - Removing others for $y_n(\mathbf{w}^T \mathbf{x}_n + b) > 1 - \xi_n$ won’t
- Cases (for SVs):
 - $\xi_n = 0$: points are $\frac{1}{\|\mathbf{w}^{\text{SVM}}\|_2}$ from the hyperplane
 - $\xi_n < 1$: correctly predicted
 - $\xi_n > 1$: wrongly predicted



Summary: Support vectors

- One of the most popular algorithms (its kernel version) before deep learning
- Theoretical guarantee on generalization error
- Convex optimization: convex objective function + **convex constraint**
 - Local optimum = global optimum
 - The solution is independent of the initialization
 - The **optimal set** is convex
 - If the objective function is **strictly convex**, then the problem has at most one optimal solution

CSE 5523: Convex Optimization + Multi-Class Classification



THE OHIO STATE UNIVERSITY

Midterm exam and HW # 3

- Midterm
 - 10/17, in class
 - Not open book, but you can bring 5 “papers” (10 pages) of cheat sheets. **Printing is fine.**
 - This means, you can practice how to make concise notes for yourself!
 - More information will be announced via Carmen!
- HW # 3
 - Due: **Due 10/15**
- Homework
 - Not meant to be a practice for your exam
 - But to verify or practice some important concepts of ML

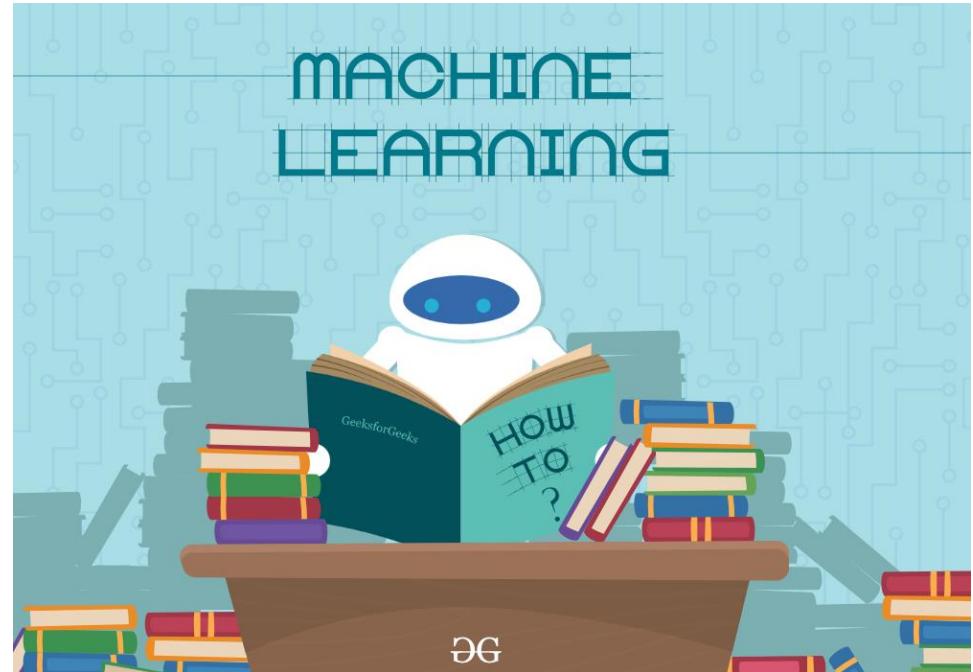
Today

Support vector machines (SVM)

- Review
- Soft-margin SVM
- Summary

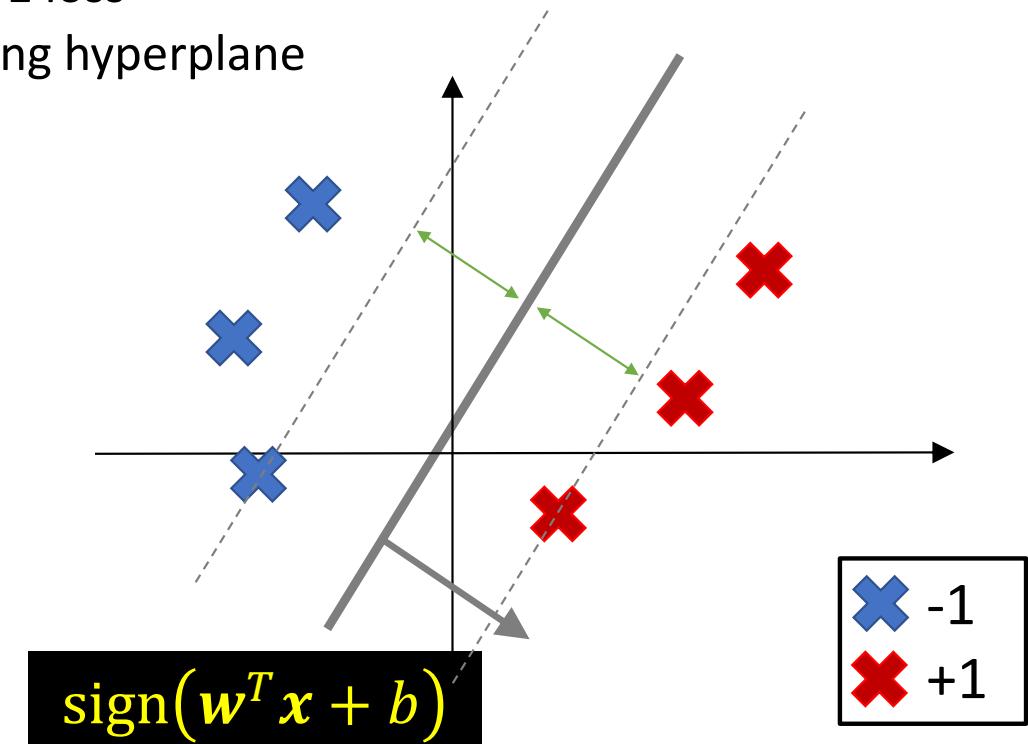
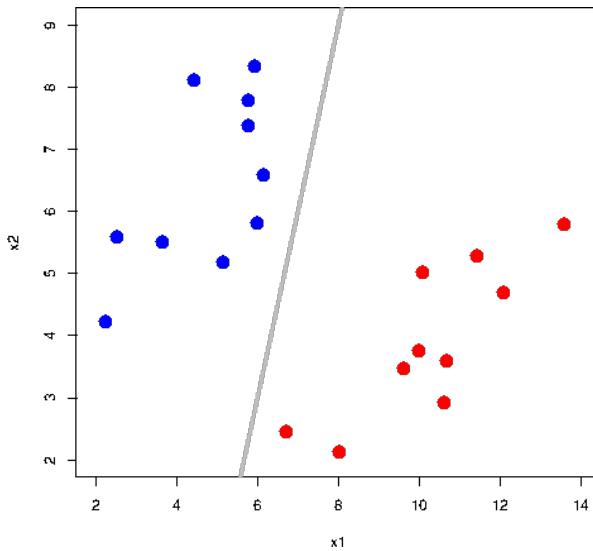
Constrained convex optimization

Multi-class classification modeling



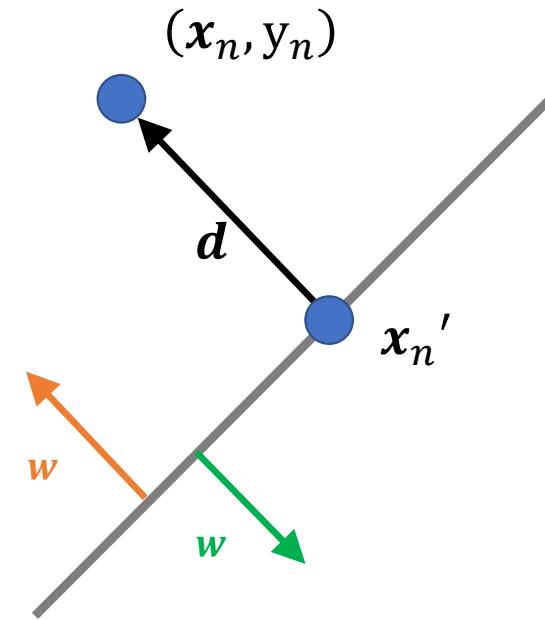
Recap: Support vector machines (SVM)

- Given a linear separable D_{tr}
 - There are infinite many (w, b) for zero 0-1 loss
 - SVM finds the maximum margin separating hyperplane



Recap: Margin

- “Sign” distance for d : $\frac{y_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|_2}$
- Properties:
 - The hyperplane and distance are positive-scale-invariant
 - $\frac{y_n(\beta \mathbf{w}^T \mathbf{x}_n + \beta b)}{\|\beta \mathbf{w}\|_2} = \frac{y_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|_2}$ for $\beta > 0$
- Margin: $\gamma = \min_{(\mathbf{x}, y) \in D_{tr}} \frac{y(\mathbf{w}^T \mathbf{x} + b)}{\|\mathbf{w}\|_2}$



$$\begin{aligned}\mathbf{w}^T \mathbf{x} + b &= 0 \\ \mathbf{w}^T \mathbf{x} + b &= 0\end{aligned}$$

Recap: Maximum margin classifier

- Training data: $D_{tr} = \{(\mathbf{x}_n \in \mathbb{R}^D, y_n \in \{+1, -1\})\}_{i=1}^N$ (linear separable)
- $(\mathbf{w}, b)^{\text{SVM}} = \operatorname{argmax}_{(\mathbf{w}, b)} \gamma(\mathbf{w}, b) = \operatorname{argmax}_{(\mathbf{w}, b)} \min_{(\mathbf{x}_n, y_n)} \frac{y_n(\mathbf{w}^T \mathbf{x}_n + b)}{\|\mathbf{w}\|_2} = \operatorname{argmax}_{(\mathbf{w}, b)} \frac{1}{\|\mathbf{w}\|_2} \min_{(\mathbf{x}_n, y_n)} y_n(\mathbf{w}^T \mathbf{x}_n + b)$
- Due to **positive-scale-invariant**, for (\mathbf{w}, b) that perfectly separate data, we can
 - constrain $\min_{(\mathbf{x}_n, y_n)} y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$
 - $(\mathbf{w}, b)^{\text{SVM}} = \operatorname{argmax}_{(\mathbf{w}, b)} \frac{1}{\|\mathbf{w}\|_2}$, s.t. $\min_{(\mathbf{x}_n, y_n)} y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1$

Recap: Maximum margin classifier

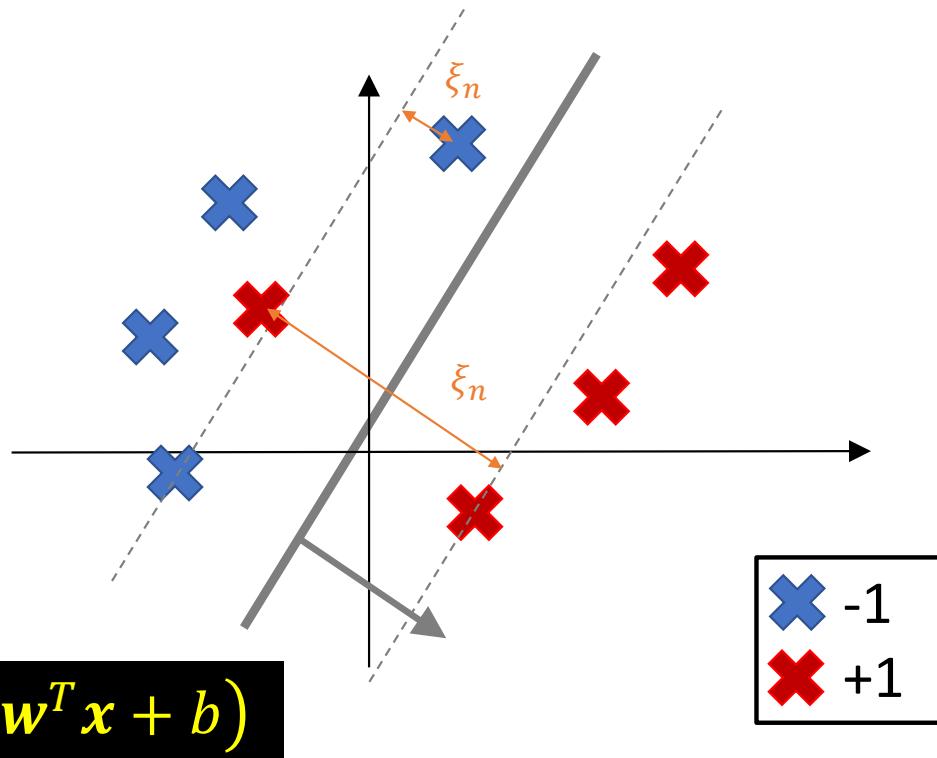
- $(\mathbf{w}, b)^{\text{SVM}} = \underset{(\mathbf{w}, b)}{\operatorname{argmax}} \frac{1}{\|\mathbf{w}\|_2} = \underset{(\mathbf{w}, b)}{\operatorname{argmin}} \|\mathbf{w}\|_2 = \underset{(\mathbf{w}, b)}{\operatorname{argmin}} \|\mathbf{w}\|_2^2 = \underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w}$
- Optimization problem:
 - $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w}, \quad \text{s.t. } \underset{(\mathbf{x}_n, y_n)}{\min} y_n (\mathbf{w}^T \mathbf{x}_n + b) = 1$
- “Equivalent” optimization problem with the same solution
 - $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w}, \quad \text{s.t. } \forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$

Recap: Soft-margin SVM

- $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w}$, s.t. $\forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$
 - No solution if the data is not linear separable (common cases)
- Allow “slight” violation with “slack” variables
- $\underset{(\mathbf{w}, b, \xi_n)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n \xi_n$, s.t. $\forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \xi_n \geq 0$
 - C is a hyper-parameter
 - How to choose it?

Recap: Slack variables

- Allow some data points to be closer to the hyperplane or classified wrongly
 - But with penalties in the objective
 - $\frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n \xi_n$
- Extreme cases:
 - If C is large, more like the hard-margin SVM
 - If C is small, favor simple models with errors



Hinge loss

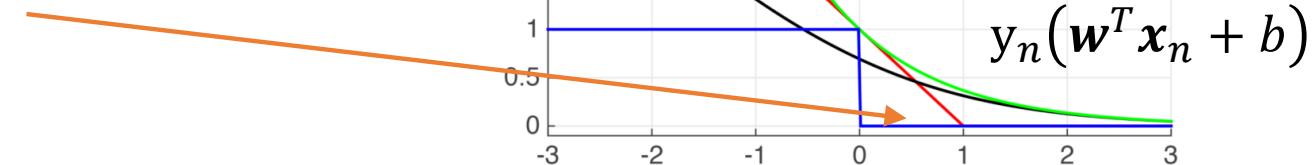
- Original soft-margin SVM

- $\underset{(\mathbf{w}, b, \xi_n)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n \xi_n$,
s.t. $\forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \xi_n \geq 0$
- Convex quadratic programming: quadratic objective + linear constraints

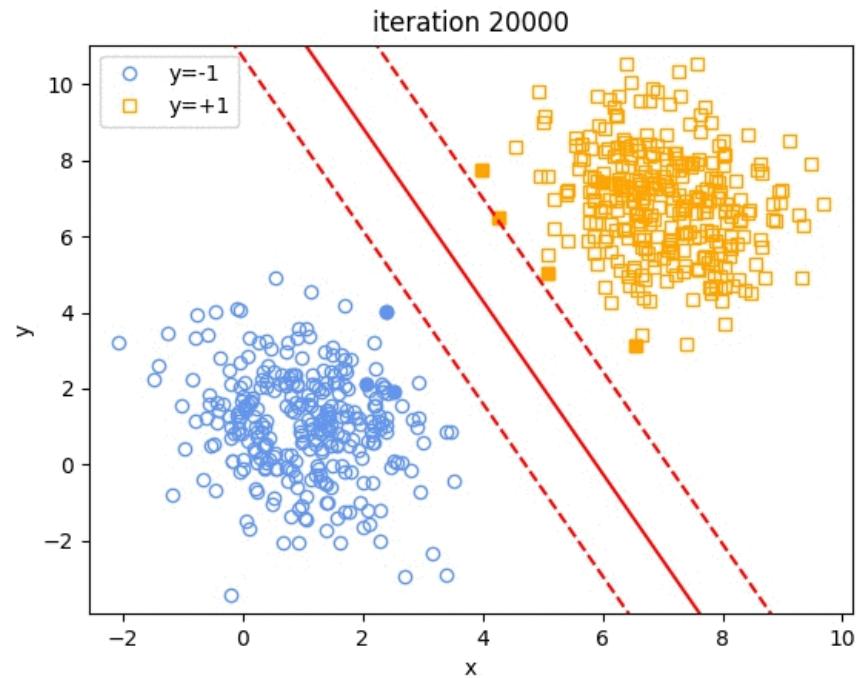
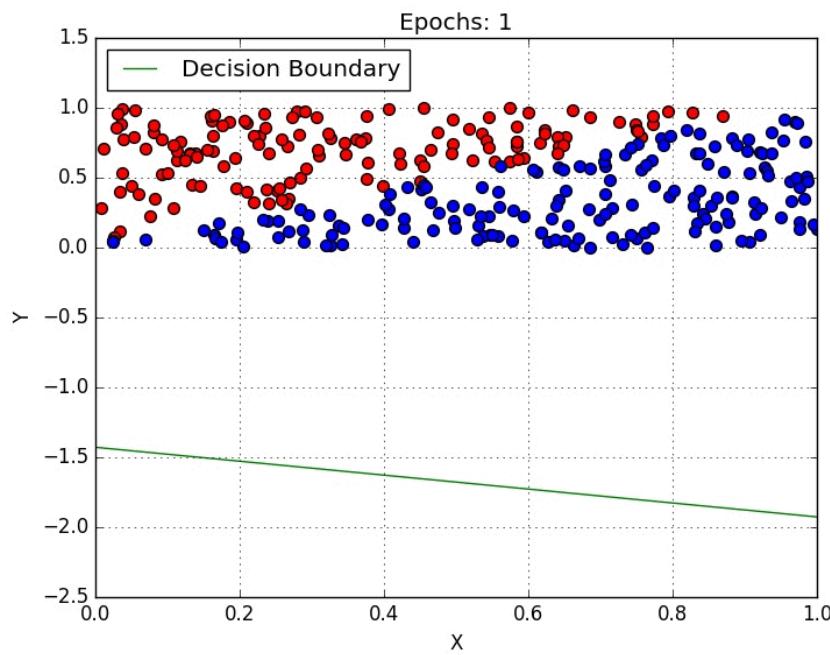
- Unconstrained formulation

- $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n [1 - y_n (\mathbf{w}^T \mathbf{x}_n + b)]_+$

- How to solve? (stochastic) gradient descent
- Penalize even for unconfident correct predictions



Iterative logistic regression and linear SVM



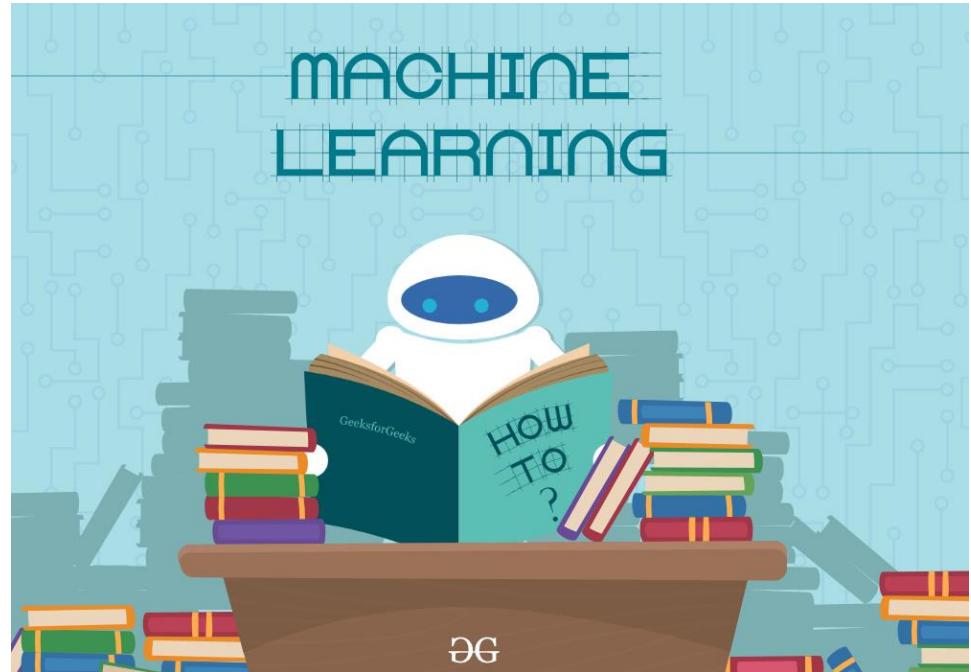
Today

Support vector machines (SVM)

- Review
- Soft-margin SVM
- Summary

Constrained convex optimization

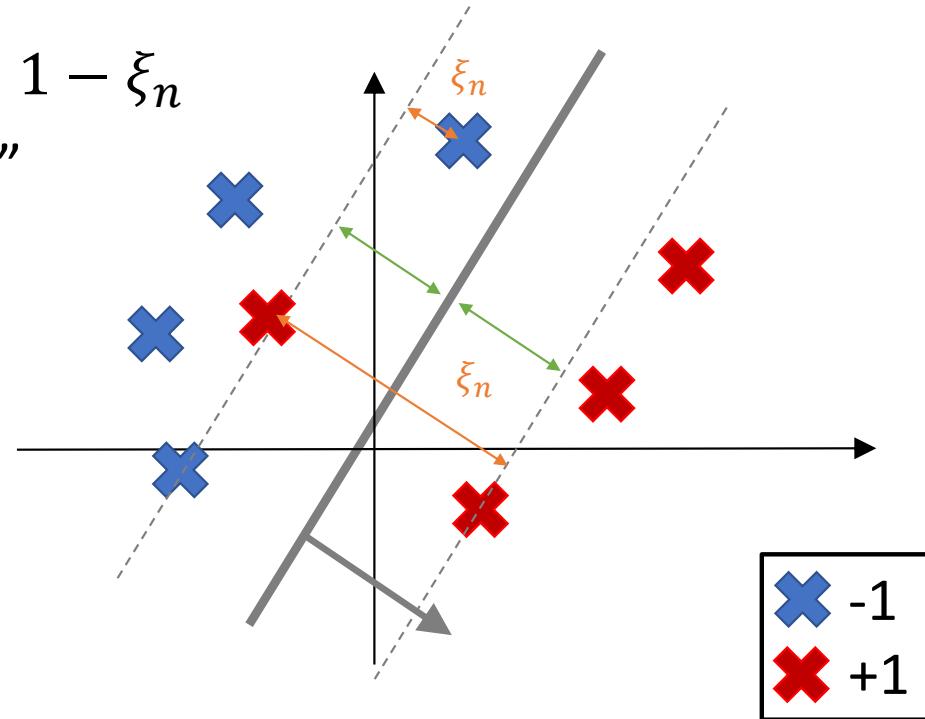
Multi-class classification modeling



Support vectors

$$\text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

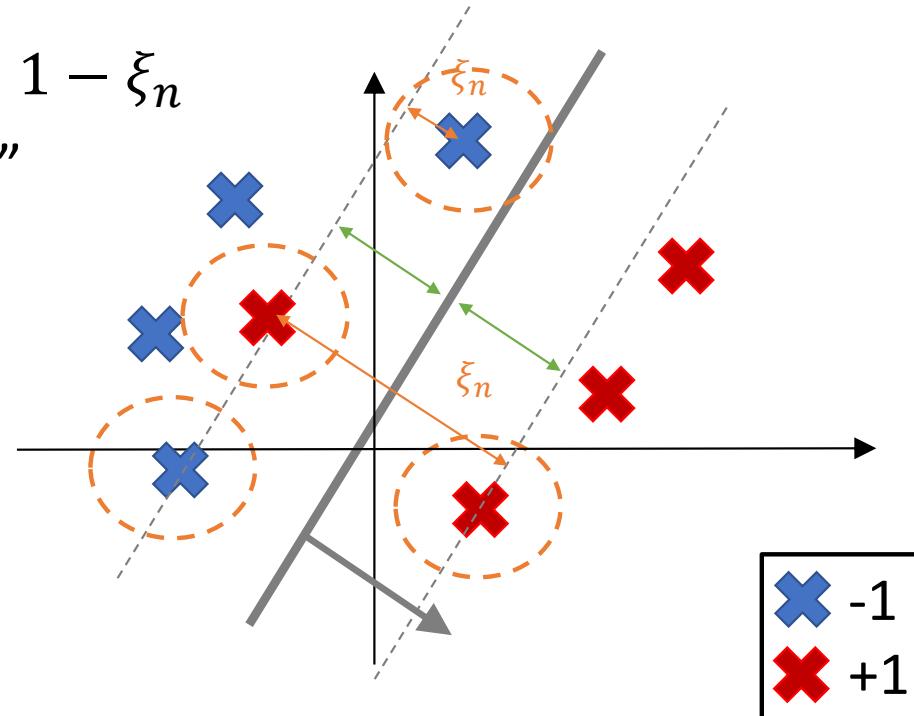
- Some (\mathbf{x}_n, y_n) will have $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1 - \xi_n$
- We refer to them as “support vectors (SV)”



Support vectors

$$\text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

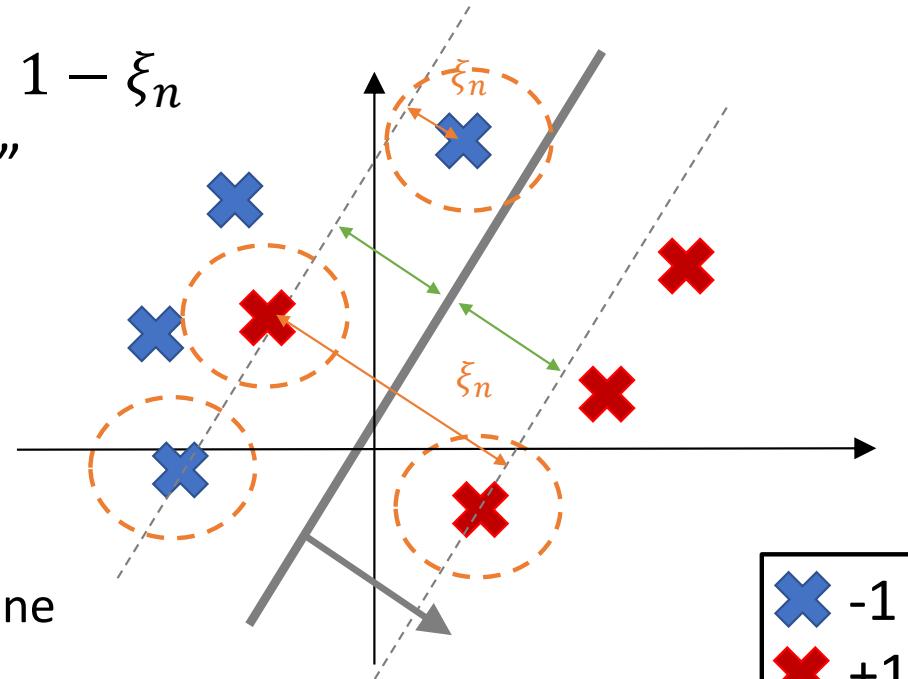
- Some (\mathbf{x}_n, y_n) will have $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1 - \xi_n$
- We refer to them as “support vectors (SV)”



Support vectors

$$\text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

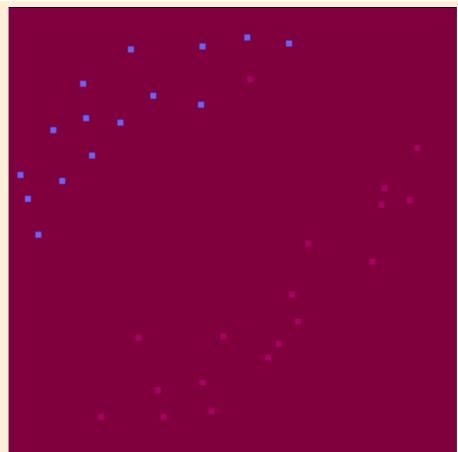
- Some (\mathbf{x}_n, y_n) will have $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1 - \xi_n$
- We refer to them as “support vectors (SV)”
- Properties: same as before
- Cases (for SVs):
 - $\xi_n = 0$: points are $\frac{1}{\|\mathbf{w}^{\text{SVM}}\|_2}$ from the hyperplane
 - $\xi_n < 1$: correctly predicted
 - $\xi_n > 1$: wrongly predicted



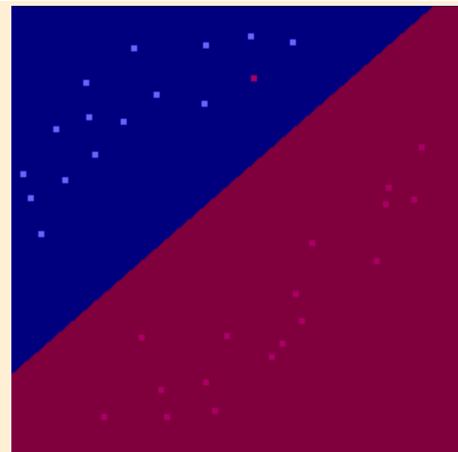
| | |
|--|----|
| | -1 |
| | +1 |

Demo

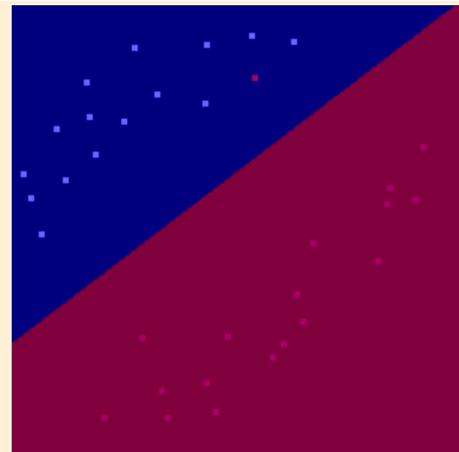
- LibSVM: <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>



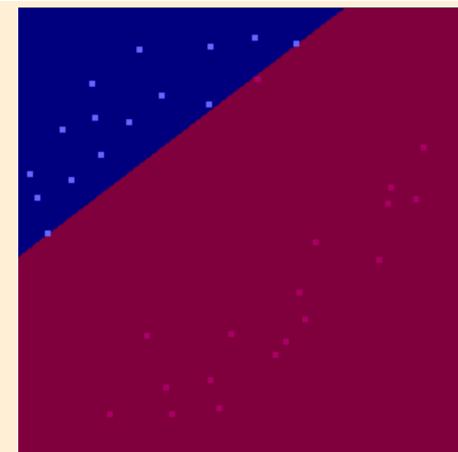
Change Run Clear -t 0 -c 0.01



Change Run Clear -t 0 -c 1



Change Run Clear -t 0 -c 100



Change Run Clear -t 0 -c 10000

$C = 0.01$

$C = 1$

$C = 100$

$C = 10000$

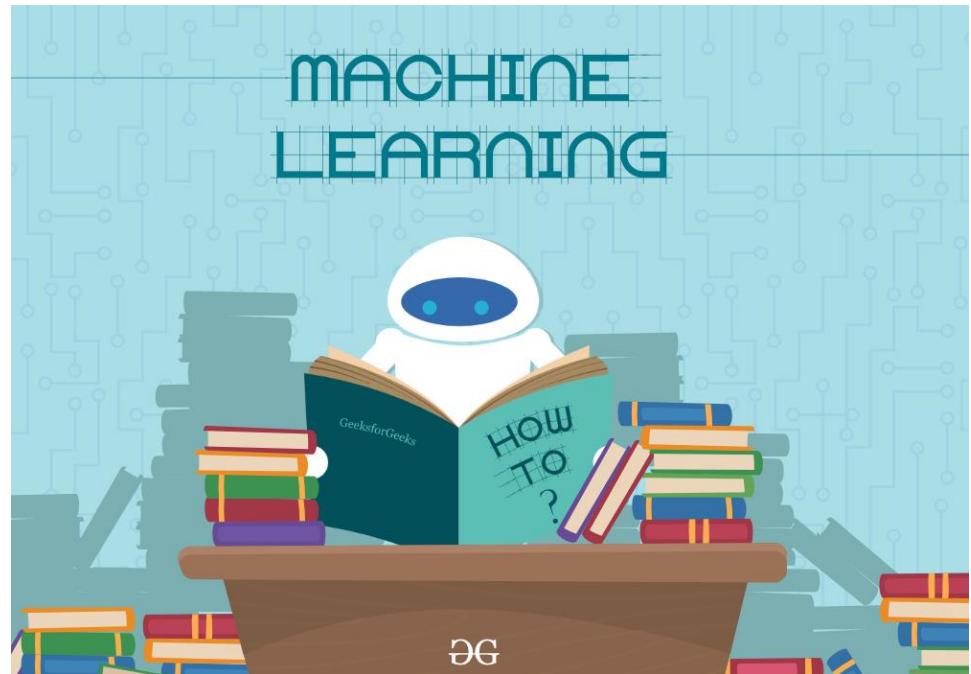
Today

Support vector machines (SVM)

- Review
- Soft-margin SVM
- Summary

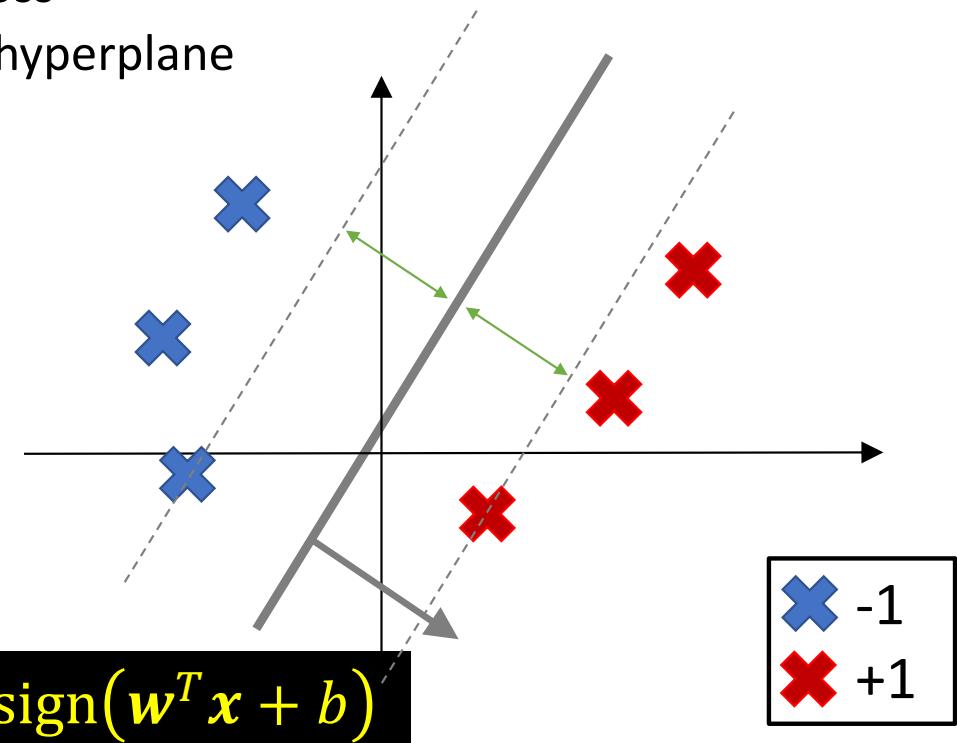
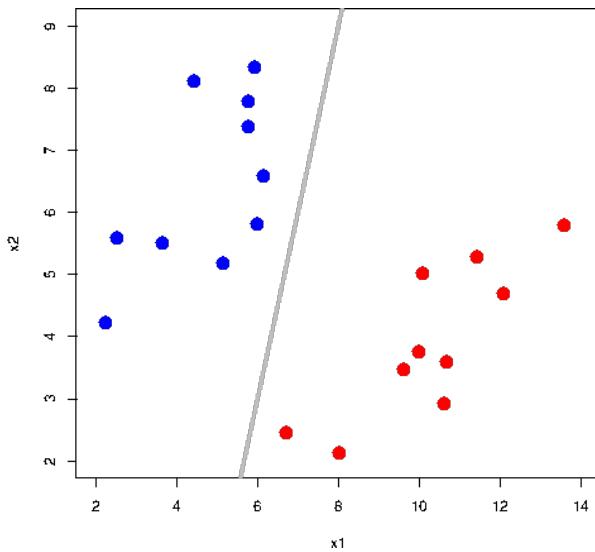
Constrained convex optimization

Multi-class classification modeling



Summary: Support vector machines (SVM)

- Given a linear separable D_{tr}
 - There are infinite many (w, b) for zero 0-1 loss
 - SVM finds the maximum margin separating hyperplane



Summary: Hard-margin vs. Soft-margin SVM

- Linear “hard-margin” SVM

- $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w}$, s.t. $y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$

- No solution if the data is not linear separable (common cases)

- Linear “soft-margin” SVM

- $\underset{(\mathbf{w}, b, \xi_n)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n \xi_n$, s.t. $\forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \xi_n \geq 0$

- Allow “slight” violation with “slack” variables ξ_n

- C is a hyper-parameter

Summary: Hinge loss

- Original soft-margin SVM

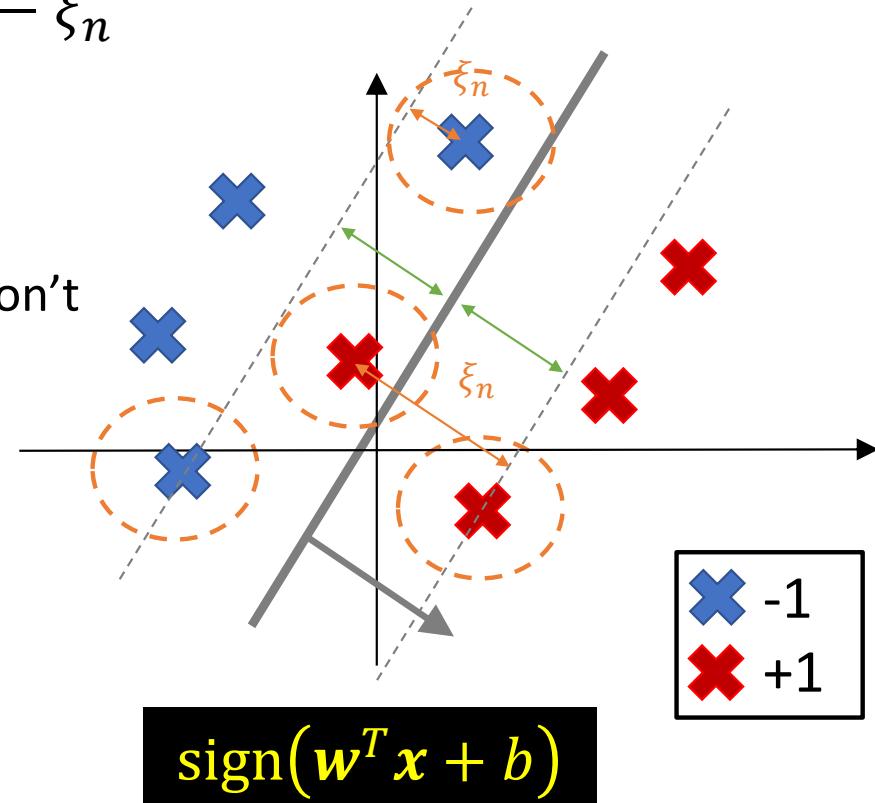
- $\underset{(\mathbf{w}, b, \xi_n)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n \xi_n,$ s.t. $\forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1 - \xi_n, \xi_n \geq 0$
- Convex quadratic programming: quadratic objective + linear constraints
- Can be solved even faster by taking the special structures of the problem into account

- Unconstrained formulation

- $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \quad \boxed{\frac{1}{2} \mathbf{w}^T \mathbf{w}} \quad + \quad \boxed{C \sum_n [1 - y_n (\mathbf{w}^T \mathbf{x}_n + b)]_+}$ training “hinge” loss regularization

Summary: Support vectors

- Some (x_n, y_n) will have $y_n(\mathbf{w}^T \mathbf{x}_n + b) = 1 - \xi_n$
 - We refer to them as “support vectors (SV)”
 - SVs define the decision boundary
 - Removing them change the decision boundary
 - Removing others for $y_n(\mathbf{w}^T \mathbf{x}_n + b) > 1 - \xi_n$ won’t
- Cases (for SVs):
 - $\xi_n = 0$: points are $\frac{1}{\|\mathbf{w}^{\text{SVM}}\|_2}$ from the hyperplane
 - $\xi_n < 1$: correctly predicted
 - $\xi_n > 1$: wrongly predicted



Summary: Support vectors

- One of the most popular algorithms (its kernel version) before deep learning
- Theoretical guarantee on generalization error
- Convex optimization: convex objective function + **convex constraint**
 - Local optimum = global optimum
 - The solution is independent of the initialization
 - The **optimal set** is convex
 - If the objective function is **strictly convex**, then the problem has at most one optimal solution

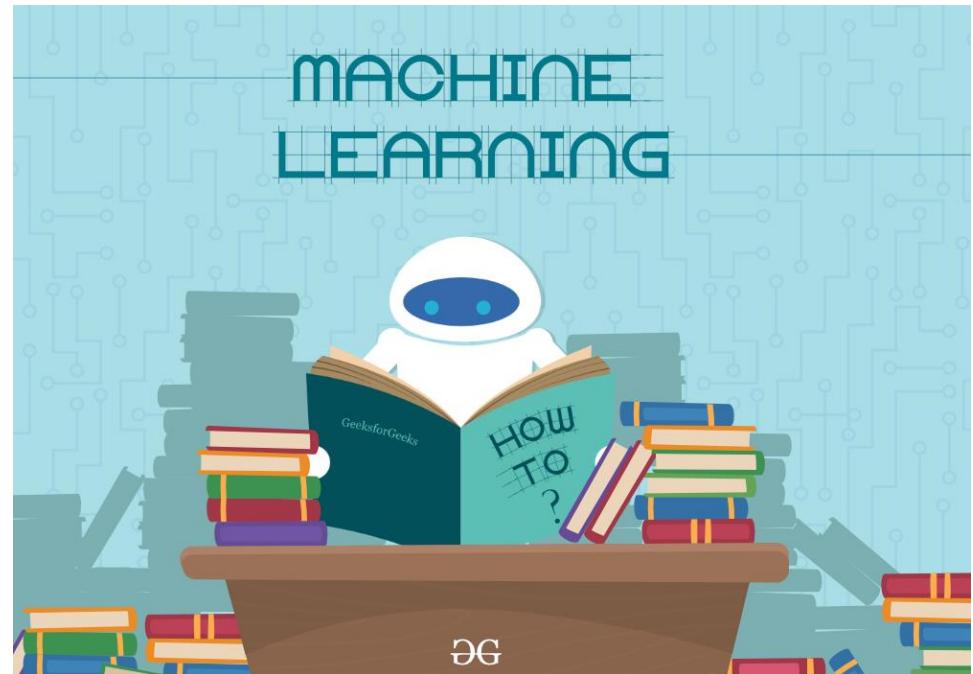
Today

Support vector machines (SVM)

- Review
- Soft-margin SVM
- Summary

Constrained convex optimization

Multi-class classification modeling



Constraint (convex) optimization

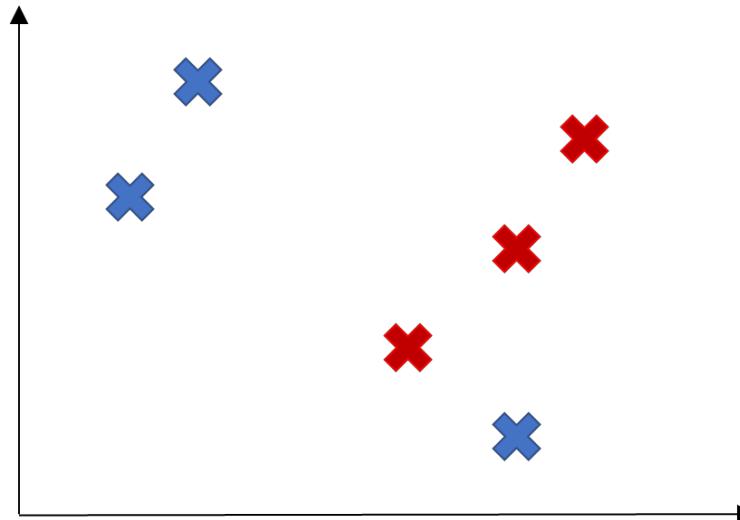
- General formulation: $\min_{\mathbf{w}} g(\mathbf{w}), \quad \text{s.t. } \mathbf{w} \in A$

Constraint (convex) optimization

- General formulation: $\min_{\mathbf{w}} g(\mathbf{w}), \quad \text{s.t. } \mathbf{w} \in \mathcal{A}$
 - \mathcal{A} is a set of feasible \mathbf{w} ; i.e., you can only choose \mathbf{w} from \mathcal{A} to minimize $g(\mathbf{w})$
 - Example: the shortest travel time from Columbus to NYC under \$50.
 - If \mathcal{A} is empty, then the problem is infeasible
 - Example: let's say there is no way that you can travel to NYC under \$50.

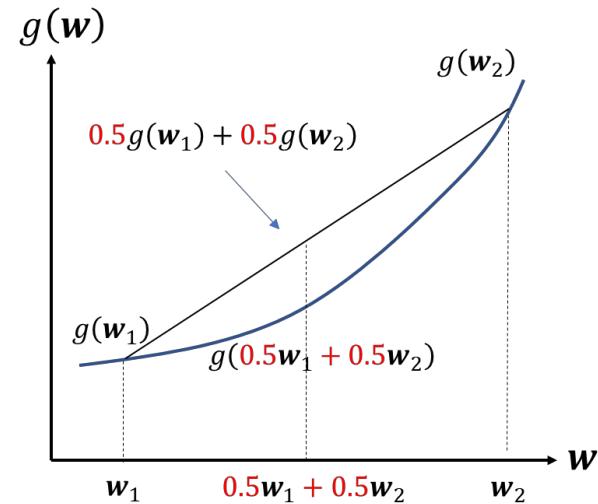
Recap: Hard-margin SVM

- $\underset{(\mathbf{w}, b)}{\operatorname{argmin}} \frac{1}{2} \mathbf{w}^T \mathbf{w}, \quad \text{s.t. } \forall n, y_n (\mathbf{w}^T \mathbf{x}_n + b) \geq 1$
 - No solution if the data is not linear separable (common cases)



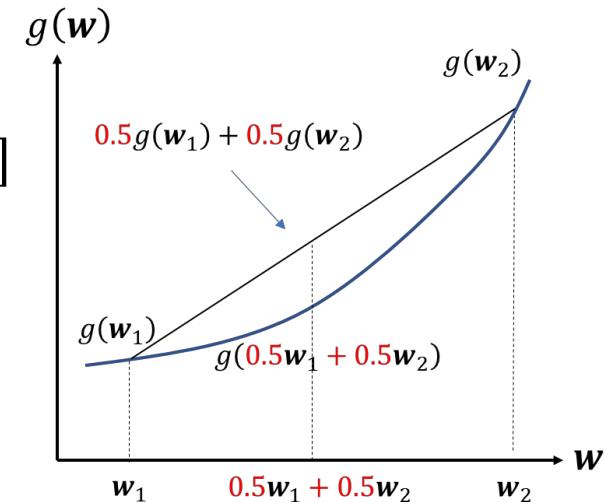
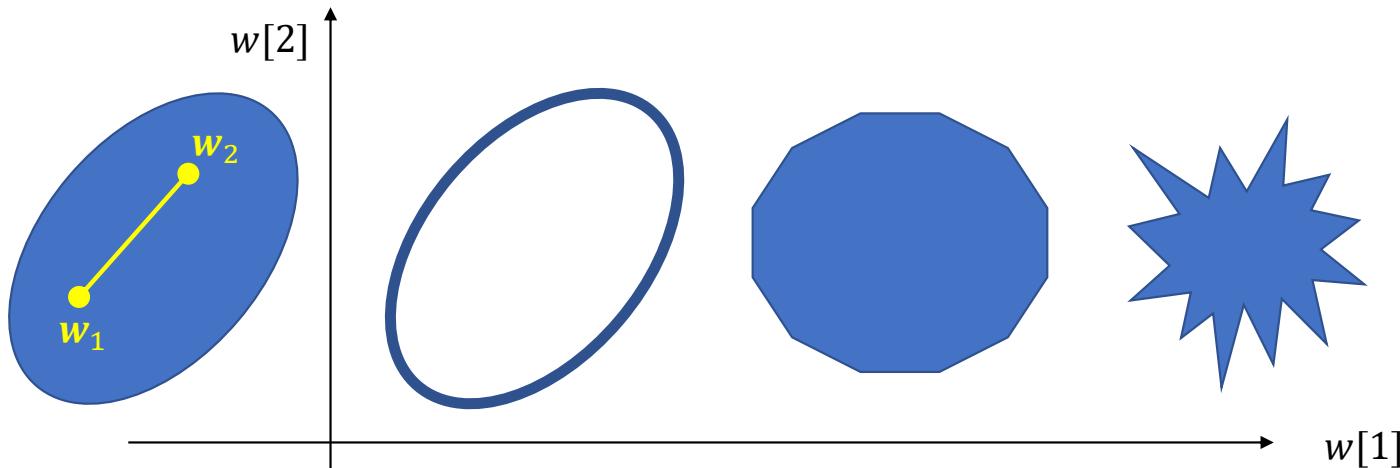
Constraint convex optimization

- $g(\mathbf{w})$ is a convex function:
 - $g(t\mathbf{w}_1 + (1 - t)\mathbf{w}_2) \leq tg(\mathbf{w}_1) + (1 - t)g(\mathbf{w}_2), \forall \mathbf{w}_1, \mathbf{w}_2; \forall t \in [0, 1]$



Constraint convex optimization

- $g(\mathbf{w})$ is a convex function:
 - $g(t\mathbf{w}_1 + (1 - t)\mathbf{w}_2) \leq tg(\mathbf{w}_1) + (1 - t)g(\mathbf{w}_2), \forall \mathbf{w}_1, \mathbf{w}_2; \forall t \in [0, 1]$
- A is a convex set:
 - If $\mathbf{w}_1 \in A$ and $\mathbf{w}_2 \in A$, then $t\mathbf{w}_1 + (1 - t)\mathbf{w}_2 \in A \forall t \in [0, 1]$

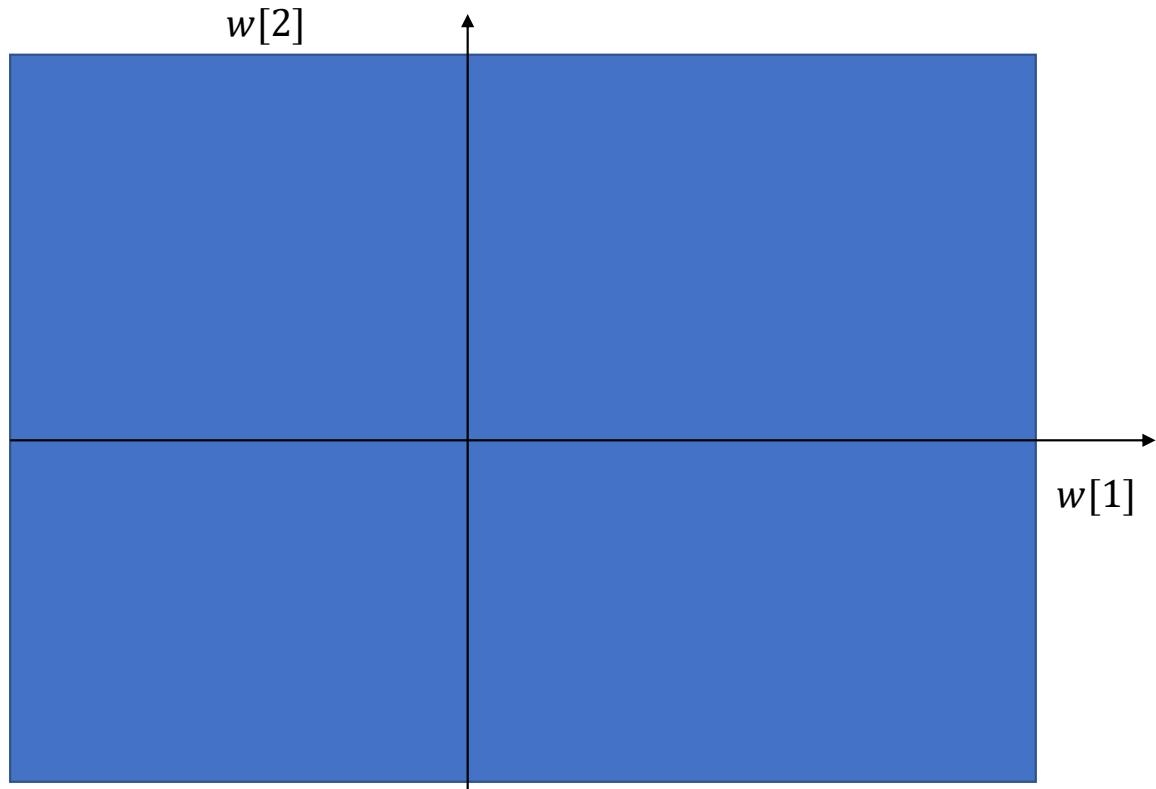


Constraint convex optimization

- The **optimal set A^*** : the set that contains the optimal solutions
 - Let $g^* = \min_w g(w)$, s.t. $w \in A$
 - $A^* = \{w \in A \text{ AND } g(w) = g^*\}$
- For convex optimization, the optimal set is convex:
 - That is, if $w_1 \in A^*$ and $w_2 \in A^*$, then $tw_1 + (1 - t)w_2 \in A^* \forall t \in [0, 1]$
- **Strictly convex** functions:
 - $g(tw_1 + (1 - t)w_2) < tg(w_1) + (1 - t)g(w_2)$, $\forall w_1, w_2; \forall t \in [0, 1]$
 - If feasible, the optimization has a unique solution

Constraint convex optimization

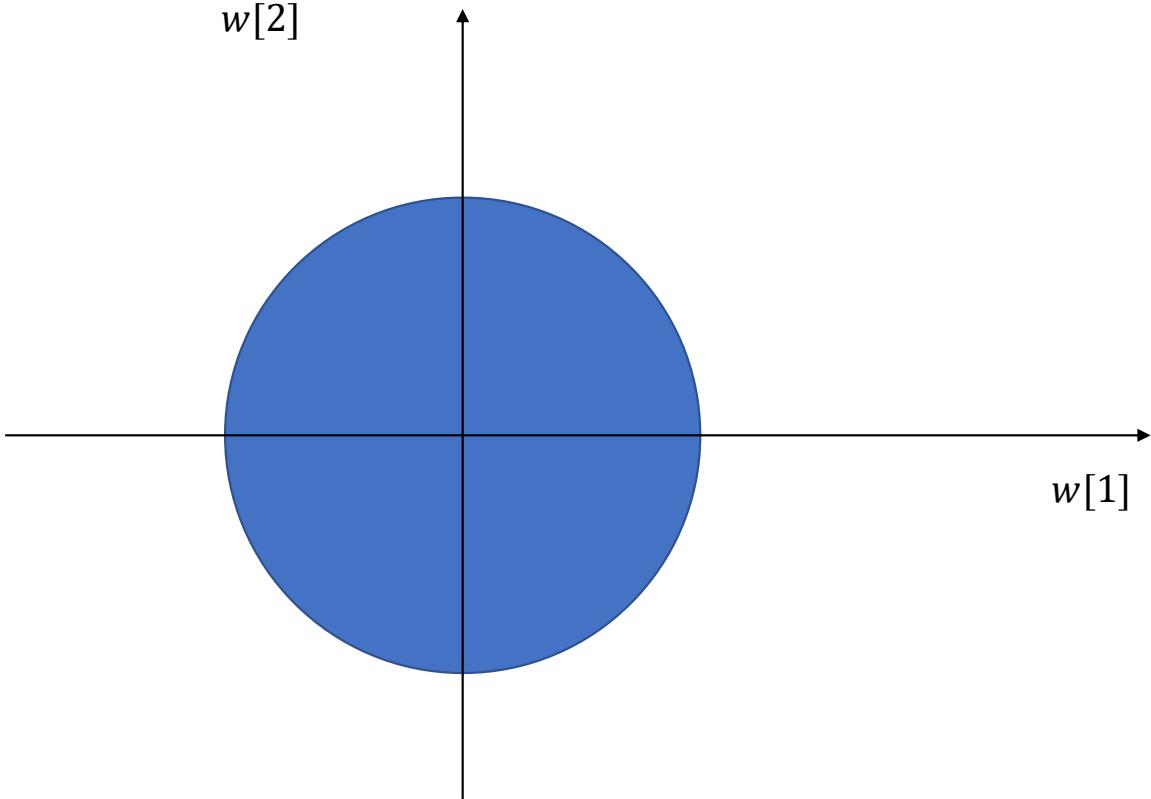
- The constraint set A is usually represented by equalities or inequalities



Constraint convex optimization

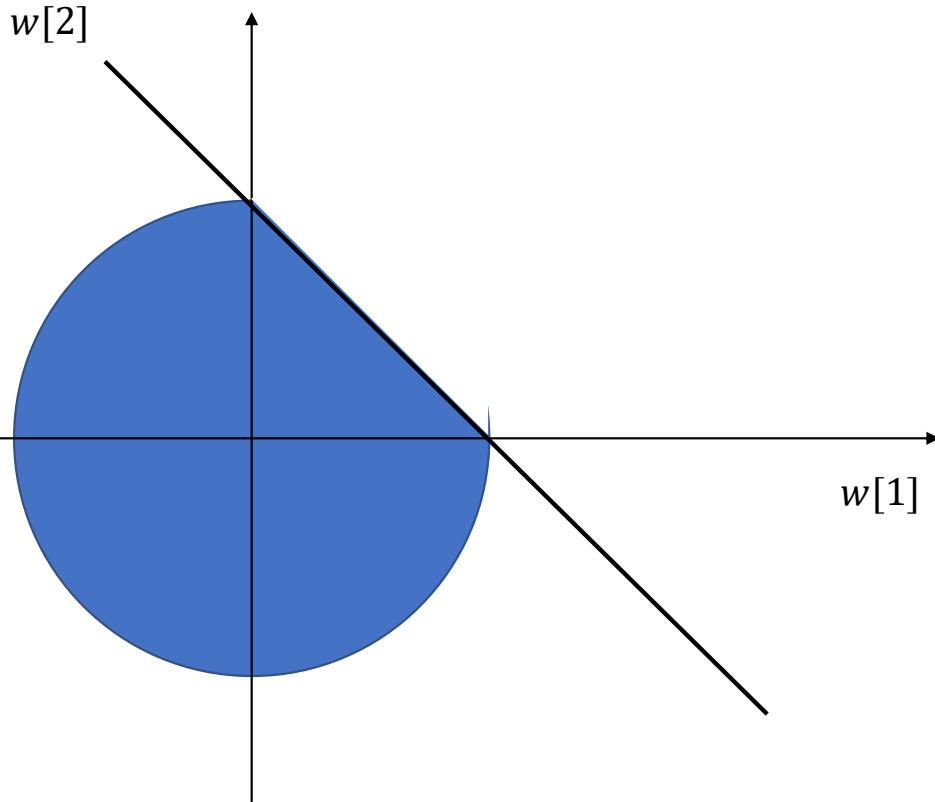
- The constraint set \mathbf{A} is usually represented by equalities or inequalities

$$w[1]^2 + w[2]^2 \leq 4$$



Constraint convex optimization

- The constraint set A is usually represented by equalities or inequalities

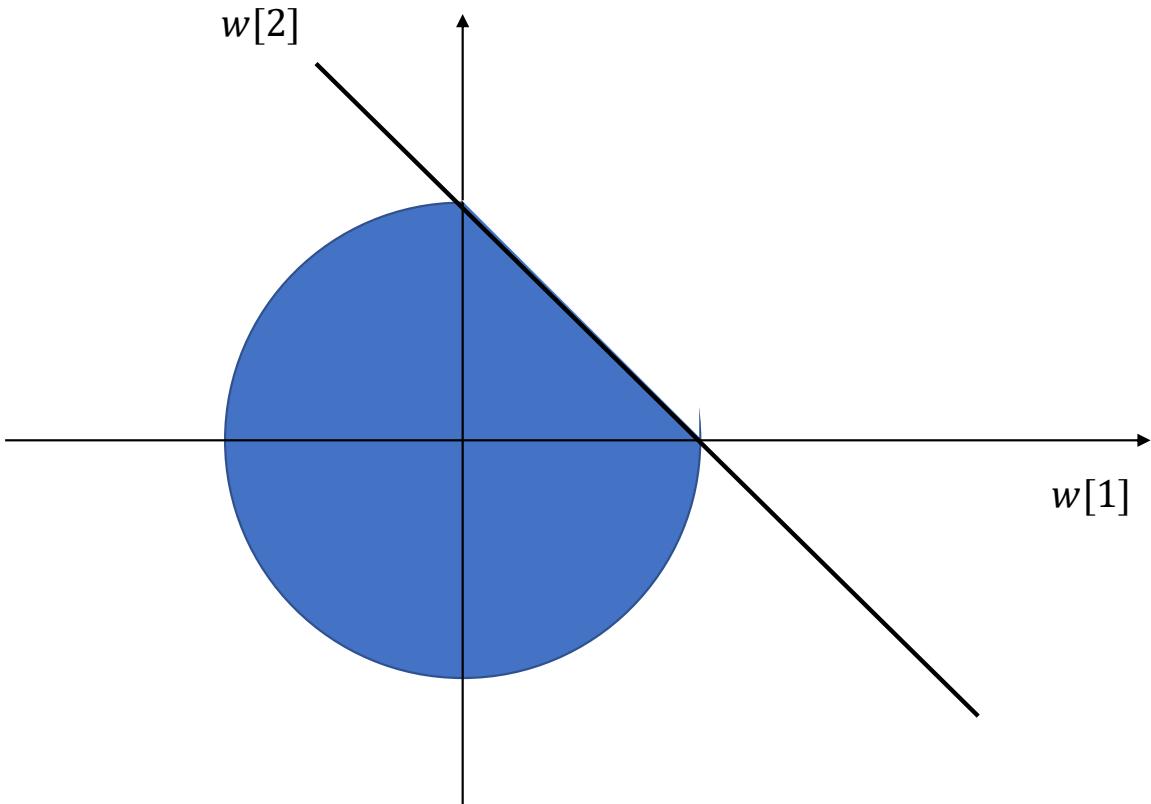


$$w[1]^2 + w[2]^2 \leq 4$$

$$2 \times w[1] + 2 \times w[2] \leq 4$$

Constraint convex optimization

- The constraint set A is usually represented by equalities or inequalities



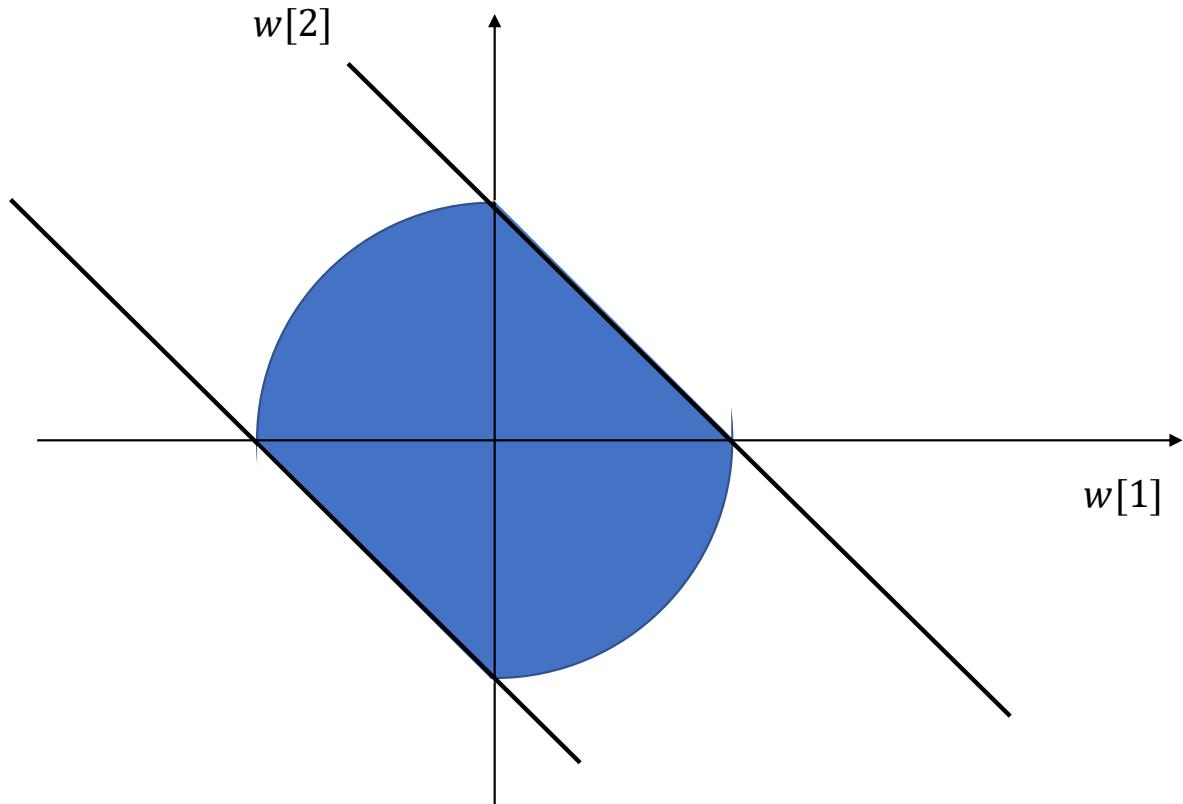
$$w[1]^2 + w[2]^2 \leq 4$$

$$2 \times w[1] + 2 \times w[2] \leq 4$$

$$2 \times w[1] + 2 \times w[2] \geq -4$$

Constraint convex optimization

- The constraint set A is usually represented by equalities or inequalities



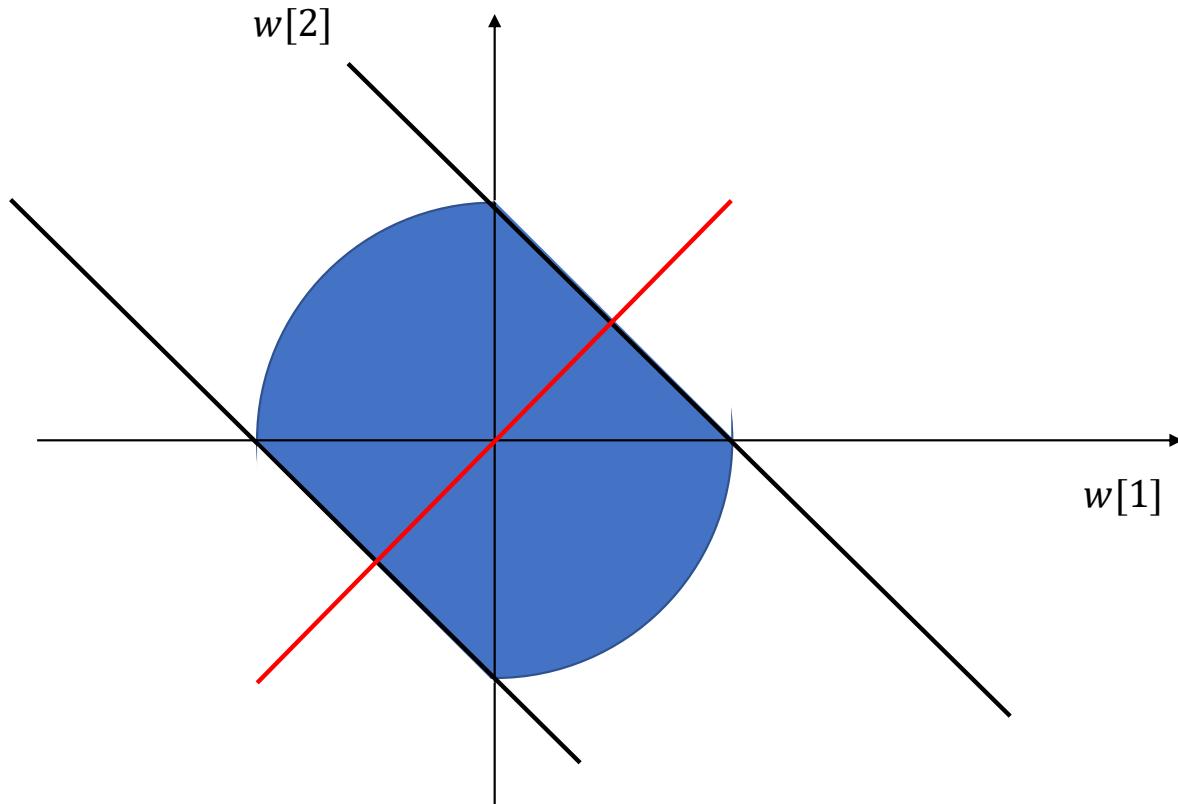
$$w[1]^2 + w[2]^2 \leq 4$$

$$2 \times w[1] + 2 \times w[2] \leq 4$$

$$2 \times w[1] + 2 \times w[2] \geq -4$$

Constraint convex optimization

- The constraint set A is usually represented by equalities or inequalities



$$w[1]^2 + w[2]^2 \leq 4$$

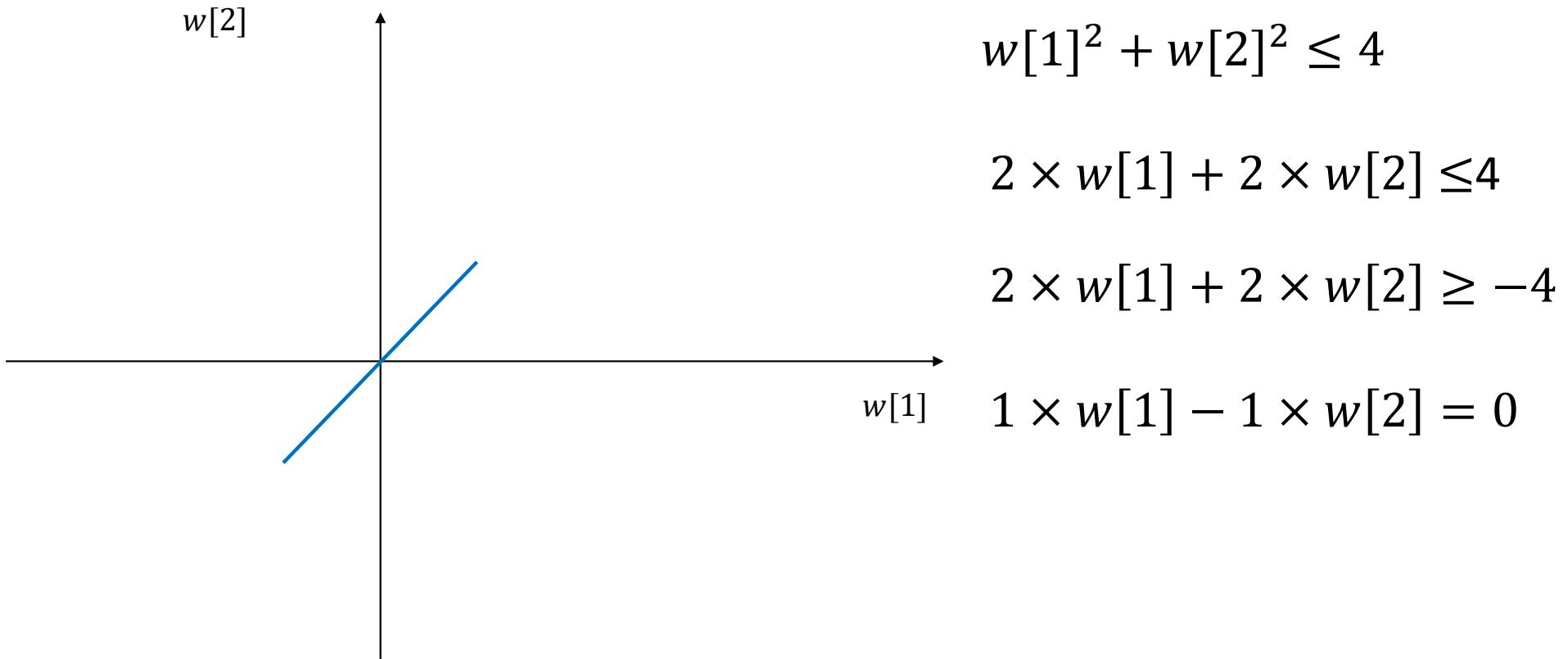
$$2 \times w[1] + 2 \times w[2] \leq 4$$

$$2 \times w[1] + 2 \times w[2] \geq -4$$

$$1 \times w[1] - 1 \times w[2] = 0$$

Constraint convex optimization

- The constraint set A is usually represented by equalities or inequalities



Constrained convex optimization

- The convex constraint set \mathbf{A} is usually represented by the “intersections of” equalities or inequalities
 - $\mathbf{A} = \{\mathbf{w} \text{ s.t. } u_i(\mathbf{w}) \leq 0, i = 1, \dots, M; \quad q_i(\mathbf{w}) = 0, i = 1, \dots, P\}$
 - $u_i(\mathbf{w})$ is a convex function
 - $q_i(\mathbf{w})$ is an affine (linear) function
 - $\mathbf{w} \in \mathbf{A}$ if it obeys all the constraints
- Example:
 - $\min_{(\mathbf{w}, b)} \frac{1}{2} \mathbf{w}^T \mathbf{w}, \quad \text{s.t. } -y_n(\mathbf{w}^T \mathbf{x}_n + b) + 1 \leq 0, \forall n$
 - $\min_{(\mathbf{w}, b, \xi_n)} \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_n \xi_n, \quad \text{s.t. } -y_n(\mathbf{w}^T \mathbf{x}_n + b) + 1 - \xi_n \leq 0, \quad -\xi_n \leq 0, \forall n$

Convex optimization

- Convex optimization: convex objective function (+ convex constrained set)
 - Local optimum = global optimum
 - The solution is independent of the initialization
 - The optimal set is convex
 - If the objective function is strictly convex, the problem has at most one optimal solution
- Example tool: <https://web.stanford.edu/~boyd/software.html>
 - If you can represent your convex optimization problem in the standard form, there are many sophisticated tool to help solve your problem.

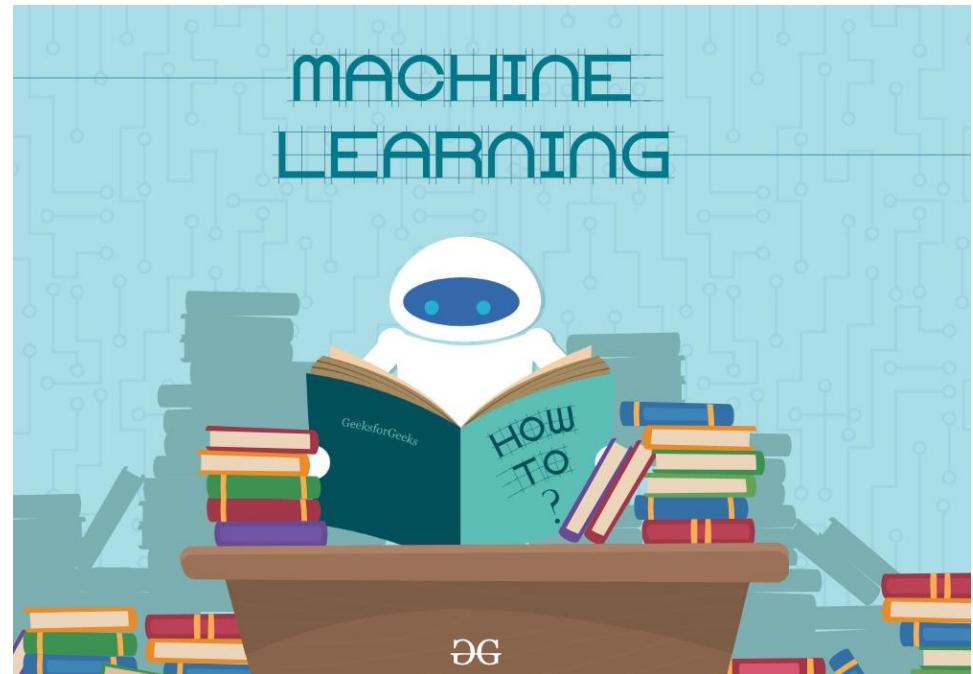
Today

Support vector machines (SVM)

- Review
- Soft-margin SVM
- Summary

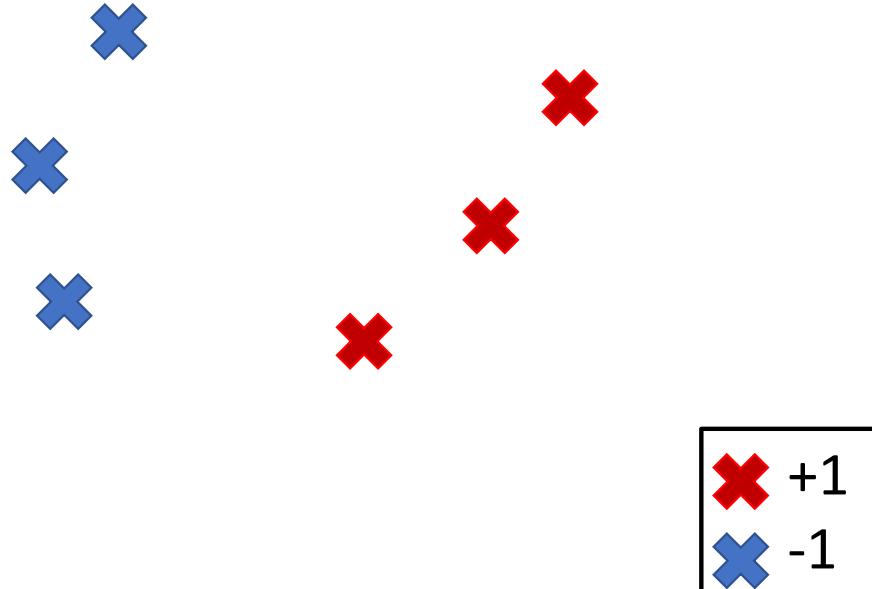
Constrained convex optimization

Multi-class classification modeling



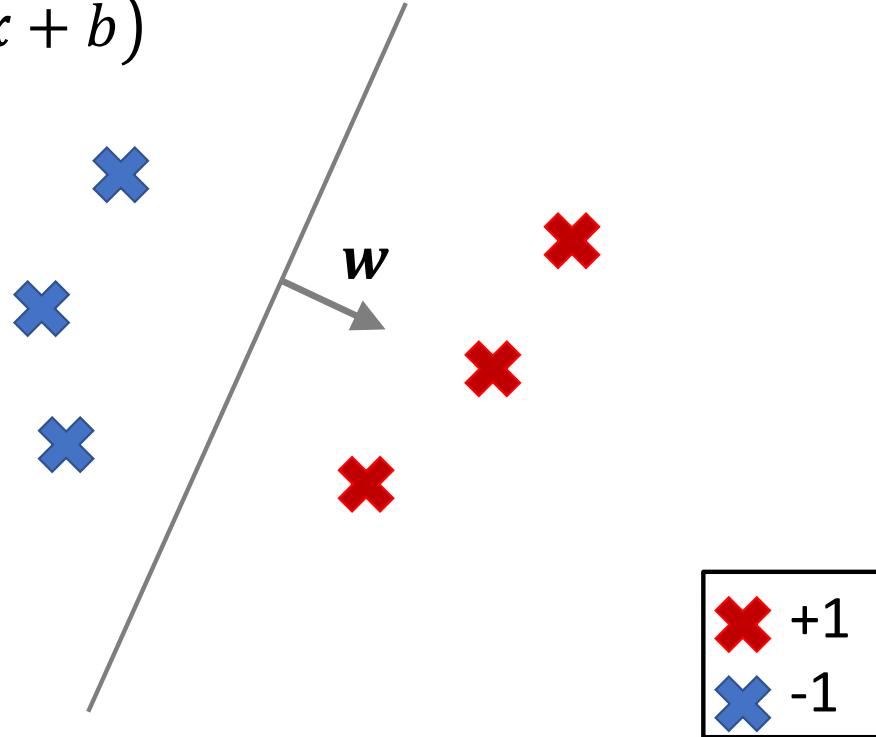
Linear classifier for binary classes

- To learn (\mathbf{w}, b)
- Decision rule: $\text{sign}(\mathbf{w}^T \mathbf{x} + b)$



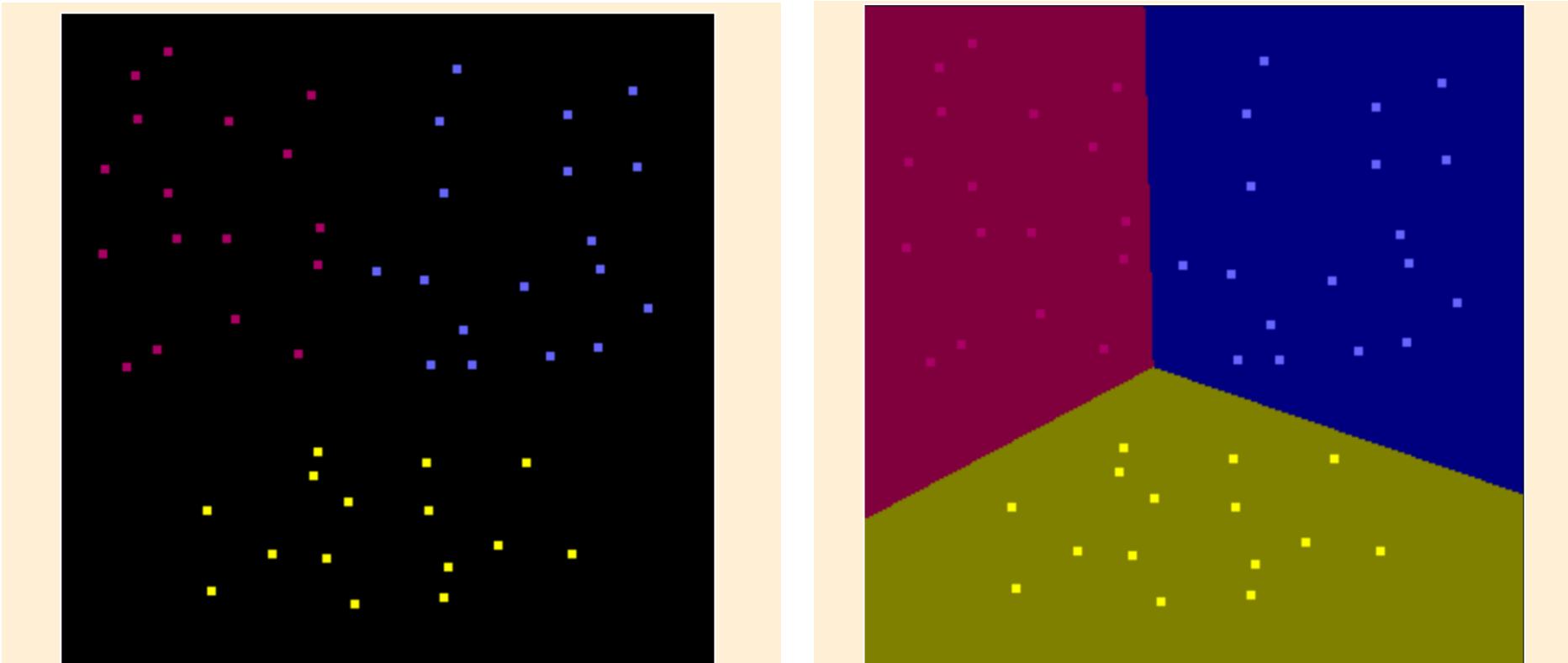
Linear classifier for binary classes

- To learn (w, b)
- Decision rule: $\text{sign}(w^T x + b)$



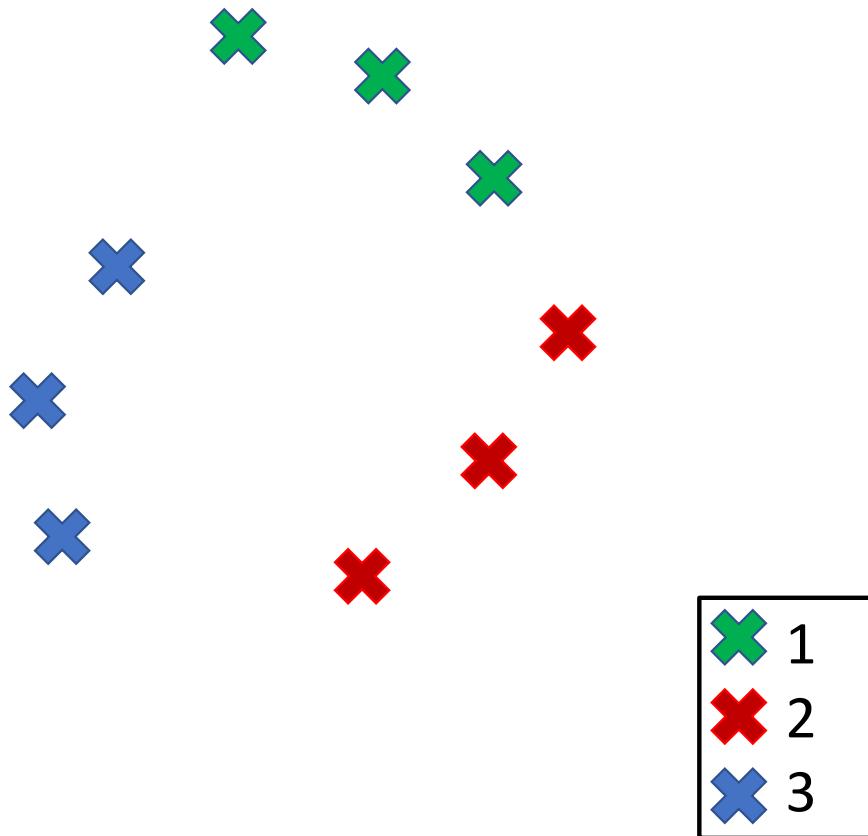
Classification

- To separate the “feature space” into class territories

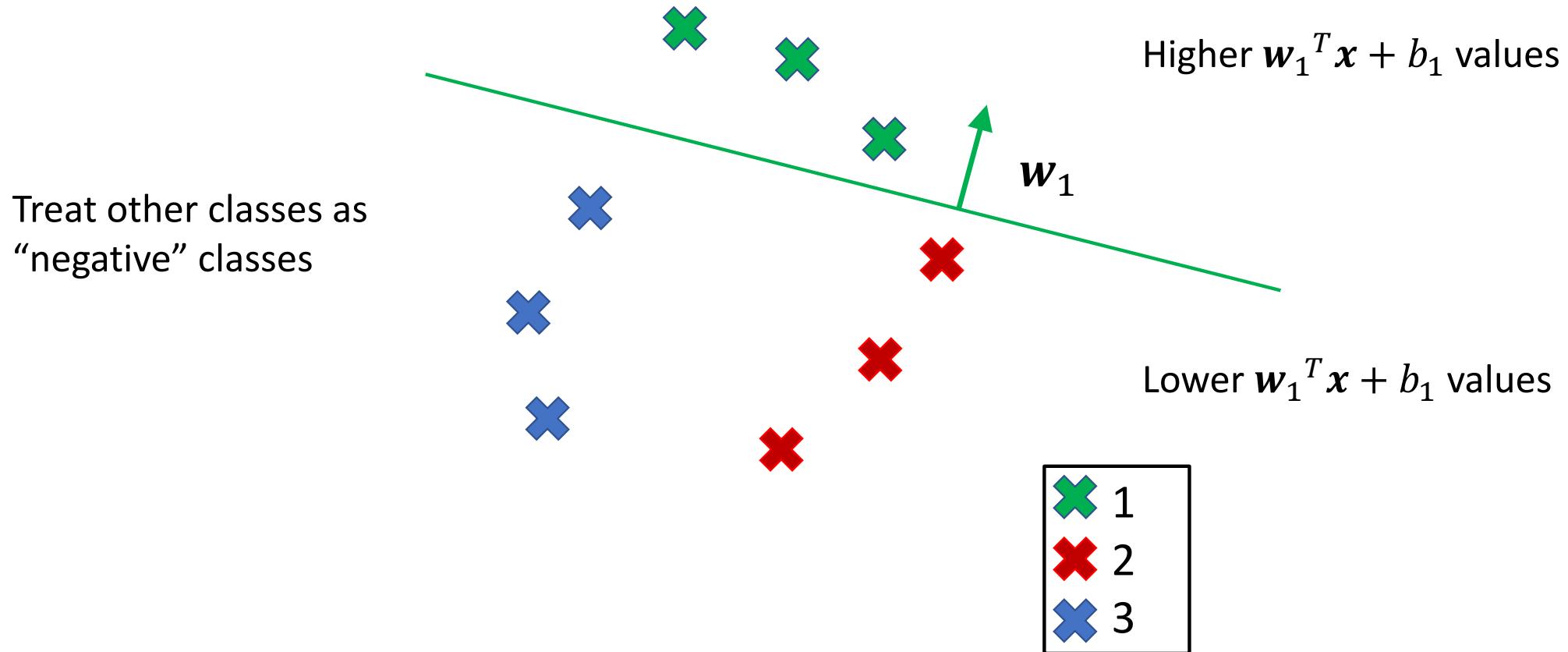


Linear classifier for multiple classes

- Linear separable?

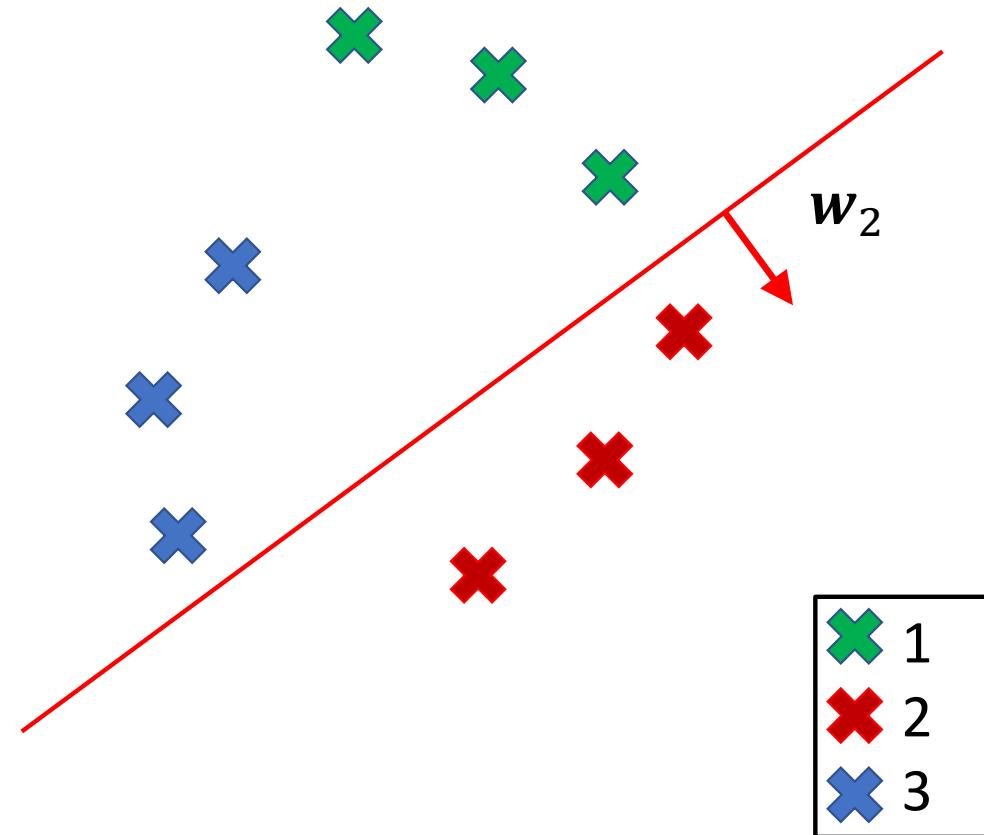


Linear classifier for multiple classes



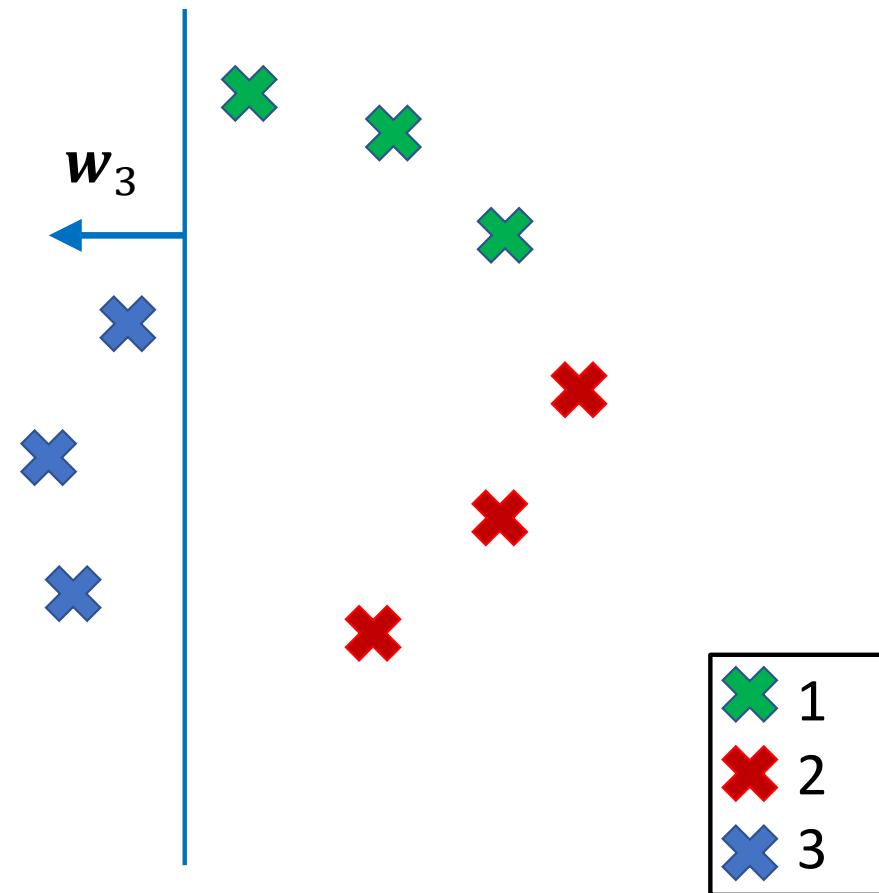
Linear classifier for multiple classes

Treat other classes as
“negative” classes



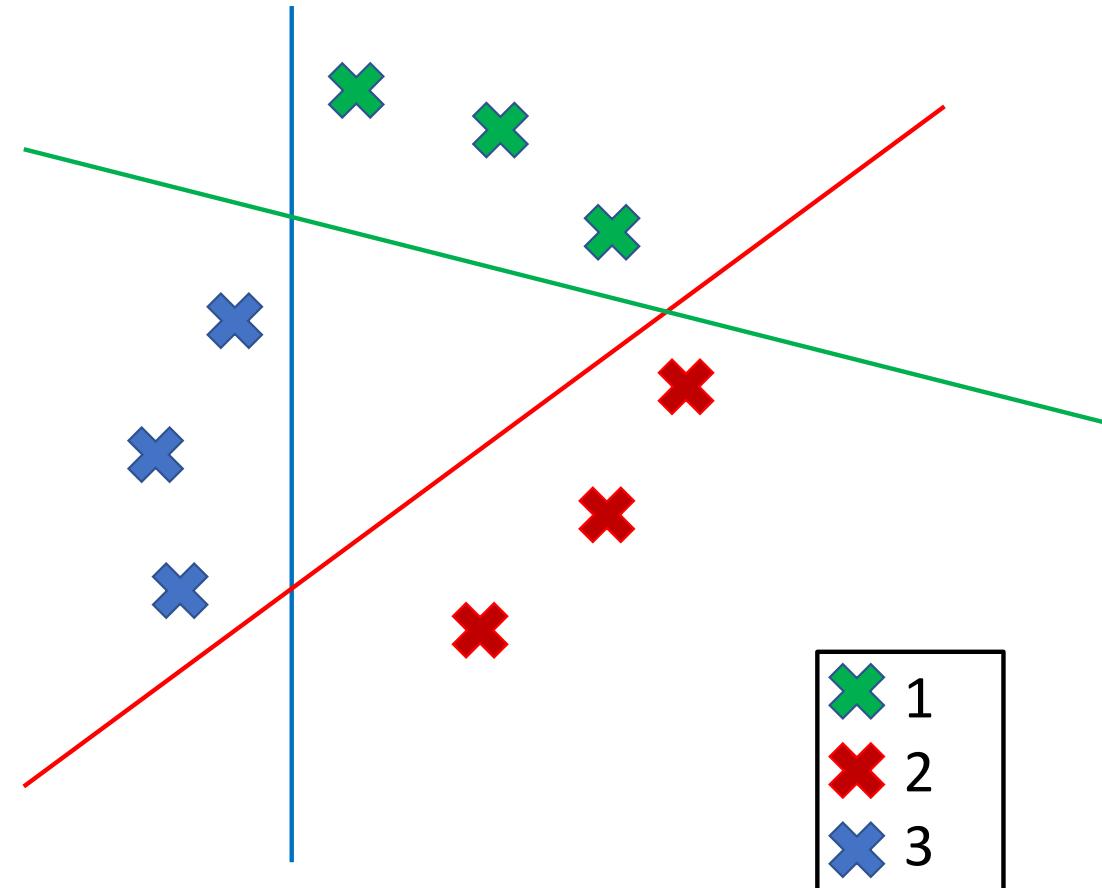
Linear classifier for multiple classes

Treat other classes as
“negative” classes



Linear classifier for multiple classes

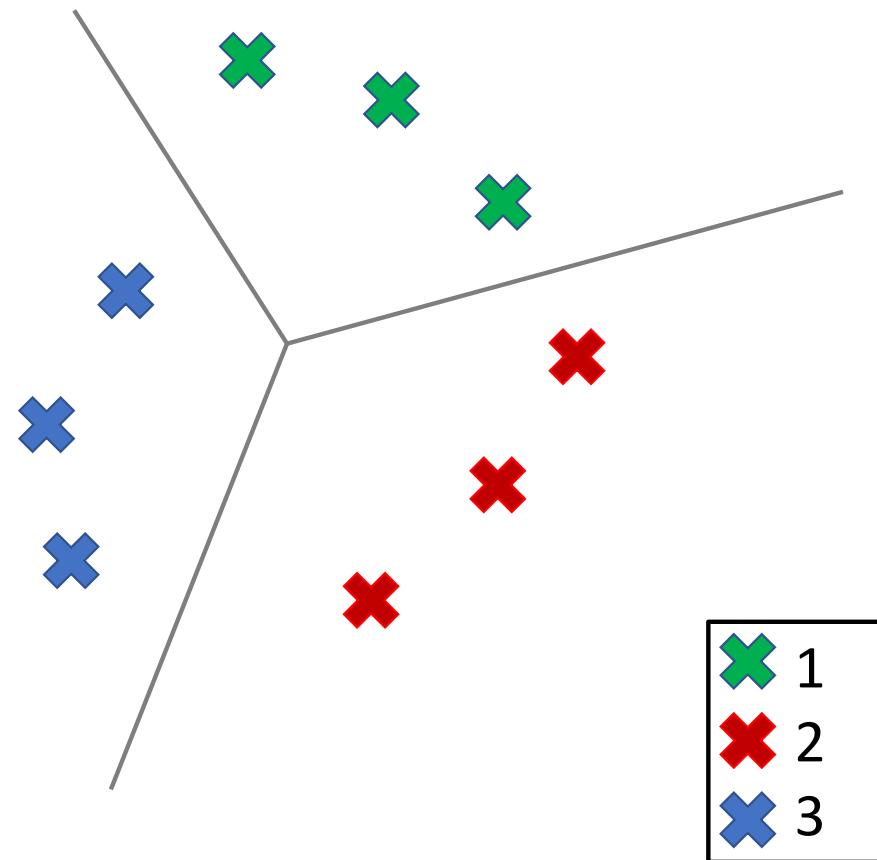
Every class c has its associated (w_c, b_c) that aims to give data from class c higher $w_c^T x + b_c$ values



Linear classifier for multiple classes

Decision rule:

$$\hat{y} = \operatorname{argmax}_c (\mathbf{w}_c^T \mathbf{x} + b_c)$$



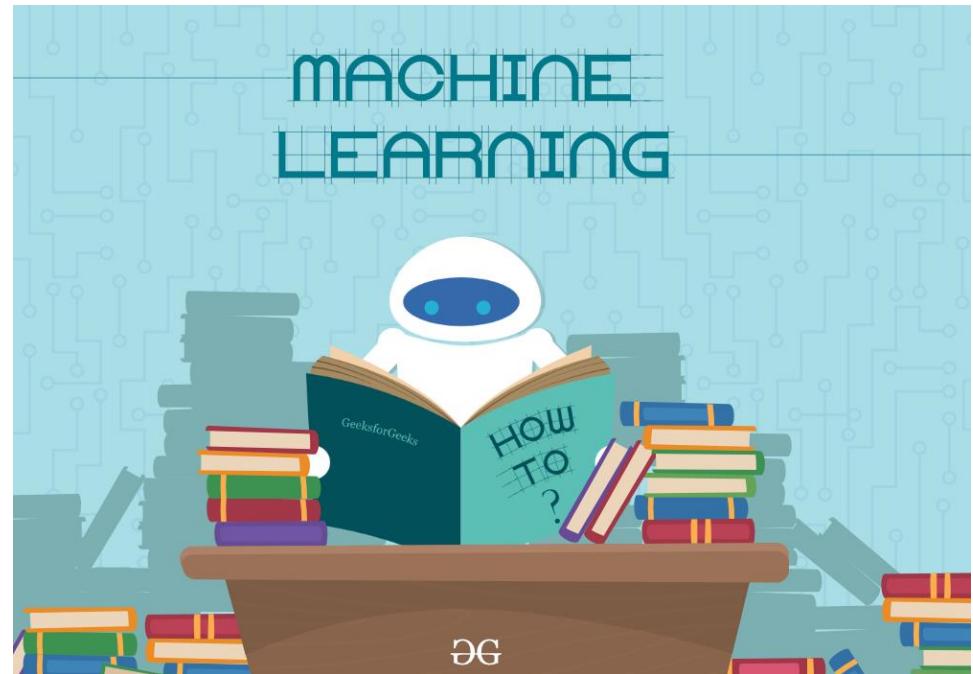
Today

Support vector machines (SVM)

- Review
- Soft-margin SVM
- Summary

Constrained convex optimization

Multi-class classification modeling



Multi-class learning

- Training data: $D_{tr} = \{ (\mathbf{x}_n \in \mathbb{R}^D, y_n) \}_{i=1}^N$
 - $y \in \{1, \dots, C\}$. (The integer C here is not the hyperparameter in SVM.)
 - $y = \mathbf{e}_c = [0, \dots, 1, \dots, 0]^T$, a C -dim **one-hot** vector with only its c -th element = 1

Multi-class learning

- Training data: $D_{tr} = \{ (\mathbf{x}_n \in \mathbb{R}^D, y_n) \}_{i=1}^N$
 - $y \in \{1, \dots, C\}$. (The integer C here is not the hyperparameter in SVM.)
 - $y = \mathbf{e}_c = [0, \dots, 1, \dots, 0]^T$, a C -dim **one-hot** vector with only its c -th element = 1
- Softmax regression
- One-vs-all (or one-vs-other) formulation: extension of SVM, logistic regression
- Crammer-Singer SVM
- Other variants: one-vs-one, DAG (Directed Acyclic Graph)

Multi-class learning

- Training data: $D_{tr} = \{ (\mathbf{x}_n \in \mathbb{R}^D, y_n) \}_{i=1}^N$
 - $y \in \{1, \dots, C\}$. (The integer C here is not the hyperparameter in SVM.)
 - $y = \mathbf{e}_c = [0, \dots, 1, \dots, 0]^T$, a C -dim **one-hot** vector with only its c -th element = 1
- Softmax regression
- One-vs-all (or one-vs-other) formulation: extension of SVM, logistic regression
- Crammer-Singer SVM
- Other variants: one-vs-one, DAG (Directed Acyclic Graph)
- **Cautions: several “notations” (like λ, e_c) may be used in different contexts**

Softmax regression (cf. logistic regression)

- **Training data:** $D_{tr} = \{ (x_n \in \mathbb{R}^D, y_n \in \{1, \dots, C\}) \}_{i=1}^N$
- **Discriminative learning:** to construct $p(Y = c|x)$

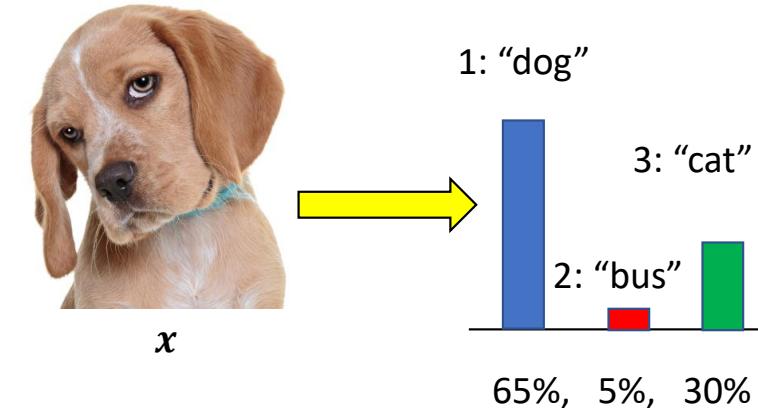
Softmax regression (cf. logistic regression)

- **Training data:** $D_{tr} = \{ (x_n \in \mathbb{R}^D, y_n \in \{1, \dots, C\}) \}_{i=1}^N$

- **Discriminative learning: to construct** $p(Y = c|x)$

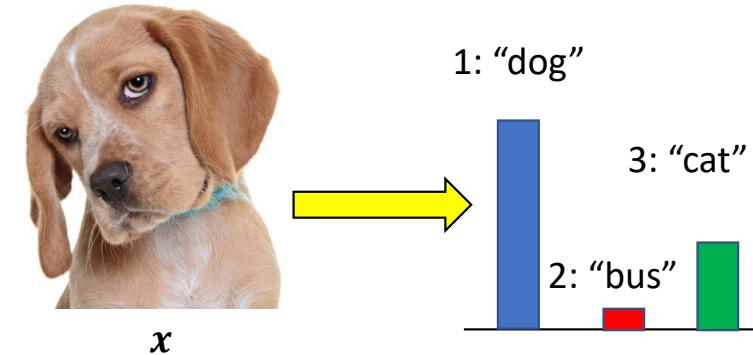
- Conditional categorical; i.e., $p(Y = c|x) = \lambda(x)[c]$

- Modeling: $\lambda(x)[c] = \frac{e^{(\mathbf{w}_c^T x + b_c)}}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T x + b_{c'})}}$



Softmax regression (cf. logistic regression)

- **Training data:** $D_{tr} = \{ (x_n \in \mathbb{R}^D, y_n \in \{1, \dots, C\}) \}_{i=1}^N$
- **Discriminative learning: to construct $p(Y = c|x)$**
 - Conditional categorical; i.e., $p(Y = c|x) = \lambda(x)[c]$
 - Modeling: $\lambda(x)[c] = \frac{e^{(\mathbf{w}_c^T \mathbf{x} + b_c)}}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x} + b_{c'})}}$
 - Properties 1: $1 \geq \lambda(x)[c] \geq 0, \sum_{c=1}^C \lambda(x)[c] = 1$
 - Properties 2: $\mathbf{w}_c^T \mathbf{x} + b_c \uparrow, p(Y = c|x) \uparrow$
 - Properties 3: $\mathbf{w}_c^T \mathbf{x} + b_c$ is usually called the decision value or logit
 - Question: what if $C = 2$? We get \mathbf{w}_1 and \mathbf{w}_2 but logistic regression only has \mathbf{w} (HW # 4)



Softmax regression

- **Learning (MLE):** $\widehat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(D_{tr}; \boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \prod_n p(y_n | x_n; \boldsymbol{\theta}) p(x_n)$

Softmax regression

- **Learning (MLE):** $\widehat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(D_{tr}; \boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \prod_n p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) p(\mathbf{x}_n)$

$$\begin{aligned} \log \prod_{n=1}^N p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) &= \sum_{n=1}^N \log p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) = \sum_{n=1}^N \log \frac{e^{(\mathbf{w}_{y_n}^T \mathbf{x}_n)}}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)}} , \text{ ignore } b_c \\ &= \sum_{n=1}^N \left[\mathbf{w}_{y_n}^T \mathbf{x}_n - \log \sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)} \right] \end{aligned}$$

Softmax regression

- **Learning (MLE):** $\widehat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(D_{tr}; \boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_n p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) p(\mathbf{x}_n)$

$$\begin{aligned} \log \prod_{n=1}^N p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) &= \sum_{n=1}^N \log p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) = \sum_{n=1}^N \log \frac{e^{(\mathbf{w}_{y_n}^T \mathbf{x}_n)}}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)}} , \text{ ignore } b_c \\ &= \sum_{n=1}^N \left[\mathbf{w}_{y_n}^T \mathbf{x}_n - \log \sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)} \right] \end{aligned}$$

- $\{\mathbf{w}_c\}_{c=1}^C \stackrel{\text{MLE}}{=} \underset{\{\mathbf{w}_c\}_{c=1}^C}{\operatorname{argmin}} \sum_{n=1}^N \left[-\mathbf{w}_{y_n}^T \mathbf{x}_n + \log \sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)} \right]$

Softmax regression: gradient

- Gradient for each data instance (\mathbf{x}_n, y_n)

$$\nabla_{\mathbf{w}_{\textcolor{red}{c}}} \left[-\mathbf{w}_{y_n}^T \mathbf{x}_n + \log \sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)} \right]$$

Softmax regression: gradient

- Gradient for each data instance (\mathbf{x}_n, y_n)

$$\begin{aligned}\nabla_{\mathbf{w}_c} \left[-\mathbf{w}_{y_n}^T \mathbf{x}_n + \log \sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)} \right] &= -\mathbf{1}[y_n == c] \times \mathbf{x}_n + \frac{\nabla_{\mathbf{w}_c} \left[\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)} \right]}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)}} \\ &= -\mathbf{1}[y_n == c] \times \mathbf{x}_n + \frac{\nabla_{\mathbf{w}_c} e^{(\mathbf{w}_c^T \mathbf{x}_n)}}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)}} = -\mathbf{1}[y_n == c] \times \mathbf{x}_n + \frac{e^{(\mathbf{w}_c^T \mathbf{x}_n)} \mathbf{x}_n}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)}} \\ &= -\mathbf{1}[y_n == c] \times \mathbf{x}_n + p(c|\mathbf{x}_n; \mathbf{w}) \mathbf{x}_n = -\{\mathbf{1}[y_n == c] - p(c|\mathbf{x}_n; \mathbf{w})\} \mathbf{x}_n\end{aligned}$$

- Note: Every (\mathbf{x}_n, y_n) contributes to the update of all $\mathbf{w}_c, \forall c \in \{1, \dots, C\}$, not only to \mathbf{w}_{y_n}

Softmax regression: cross entropy loss

- **Learning (MLE):** $\widehat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(D_{tr}; \boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \prod_n p(y_n | x_n; \boldsymbol{\theta}) p(x_n)$

$$-\log \prod_{n=1}^N p(y_n | x_n; \boldsymbol{\theta}) = -\sum_{n=1}^N \log p(y_n | x_n; \boldsymbol{\theta})$$

Softmax regression: cross entropy loss

- **Learning (MLE):** $\widehat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(D_{tr}; \boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \prod_n p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) p(\mathbf{x}_n)$

$$\begin{aligned} -\log \prod_{n=1}^N p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) &= -\sum_{n=1}^N \log p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) \\ &= \sum_{n=1}^N \sum_{c=1}^C -\mathbf{1}[y_n == c] \times \log p(c | \mathbf{x}_n; \boldsymbol{\theta}) \end{aligned}$$

Softmax regression: cross entropy loss

- **Learning (MLE):** $\widehat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(D_{tr}; \boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_n p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) p(\mathbf{x}_n)$

$$\begin{aligned} -\log \prod_{n=1}^N p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) &= -\sum_{n=1}^N \log p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) \\ &= \sum_{n=1}^N \sum_{c=1}^C -\mathbf{1}[y_n == c] \times \log p(c | \mathbf{x}_n; \boldsymbol{\theta}) \\ &= \sum_{n=1}^N \sum_{c=1}^C -\mathbf{e}_{y_n}[c] \times \log p(c | \mathbf{x}_n; \boldsymbol{\theta}) = -\sum_{n=1}^N \mathbf{e}_{y_n}^T \begin{bmatrix} \log p(1 | \mathbf{x}_n; \boldsymbol{\theta}) \\ \vdots \\ \log p(C | \mathbf{x}_n; \boldsymbol{\theta}) \end{bmatrix} \end{aligned}$$

Softmax regression: cross entropy loss

- Learning (MLE): $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} p(D_{tr}; \theta) = \underset{\theta}{\operatorname{argmax}} \prod_n p(y_n | x_n; \theta) p(x_n)$

$$-\log \prod_{n=1}^N p(y_n | x_n; \theta) = -\sum_{n=1}^N \log p(y_n | x_n; \theta)$$

$$= \sum_{n=1}^N \sum_{c=1}^C -\mathbf{1}[y_n == c] \times \log p(c | x_n; \theta)$$

$$= \sum_{n=1}^N \sum_{c=1}^C -\boxed{e_{y_n}[c]} \times \boxed{\log p(c | x_n; \theta)} = -\sum_{n=1}^N e_{y_n}^T \begin{bmatrix} \log p(1 | x_n; \theta) \\ \vdots \\ \log p(C | x_n; \theta) \end{bmatrix}$$

target prediction

Softmax regression: cross entropy loss

- Learning (MLE): $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} p(D_{tr}; \theta) = \underset{\theta}{\operatorname{argmax}} \prod_n p(y_n | x_n; \theta) p(x_n)$

$$\begin{aligned}& -\log \prod_{n=1}^N p(y_n | x_n; \theta) = -\sum_{n=1}^N \log p(y_n | x_n; \theta) \\&= \sum_{n=1}^N \sum_{c=1}^C -\mathbf{1}[y_n == c] \times \log p(c | x_n; \theta) \\&= \sum_{n=1}^N \sum_{c=1}^C -\boxed{\mathbf{e}_{y_n}[c]} \times \boxed{\log p(c | x_n; \theta)} = -\sum_{n=1}^N \mathbf{e}_{y_n}^T \begin{bmatrix} \log p(1 | x_n; \theta) \\ \vdots \\ \log p(C | x_n; \theta) \end{bmatrix}\end{aligned}$$

target prediction

- Re-write the objective
- Want the target and the prediction to be similar (i.e., to have high inner product)

Softmax regression: prediction

- Prediction by Bayes optimal classifier

$$\circ \hat{y} = \operatorname{argmax}_c p(c|x; \theta) = \operatorname{argmax}_c \frac{e^{(\mathbf{w}_{\textcolor{red}{c}}^T \mathbf{x})}}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x})}} = \operatorname{argmax}_c e^{(\mathbf{w}_{\textcolor{red}{c}}^T \mathbf{x})} = \operatorname{argmax}_c \mathbf{w}_{\textcolor{red}{c}}^T \mathbf{x}$$

Midterm exam and HW # 3

- Midterm
 - 10/17, in class
 - Not open book, but you can bring 5 “papers” (10 pages) of cheat sheets. **Printing is fine.**
 - This means, you can practice how to make concise notes for yourself!
 - **More information will be announced via Carmen!**
- HW # 3
 - Due (problem): **10/15**
 - Due (programming): **10/22**
- Homework
 - Not meant to be a practice for your exam
 - But to verify or practice some important concepts of ML

HW regrading

- TA will create Google forms about it.

CSE 5523: Multi-class + ERM



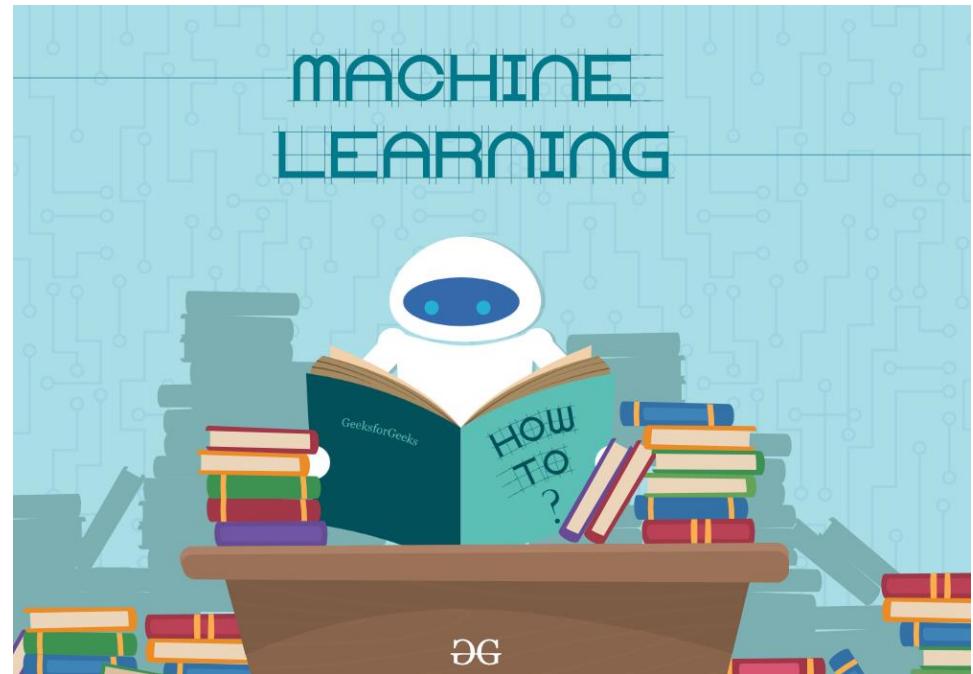
THE OHIO STATE UNIVERSITY

Today

Multi-class classification

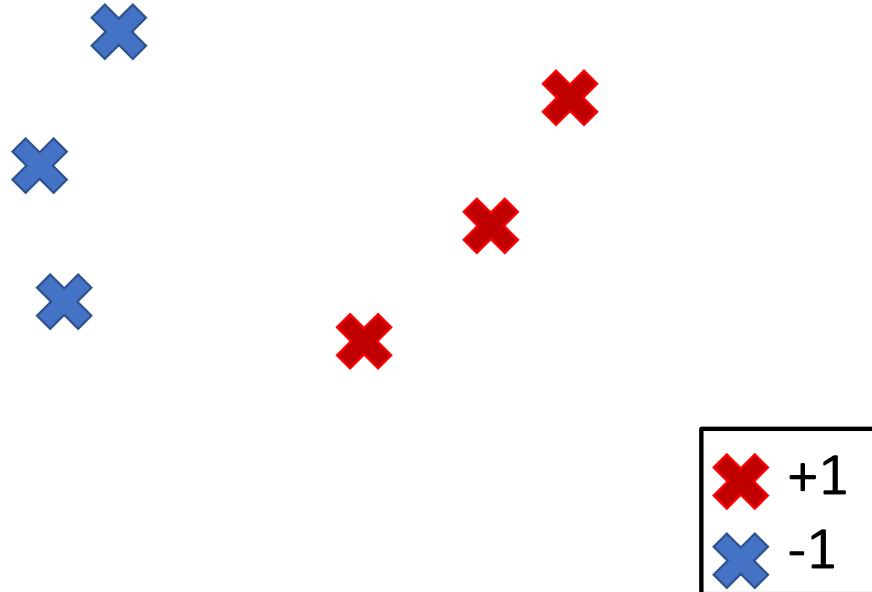
- Basics
- Modeling + Learning
- Extension

ERM + regularization



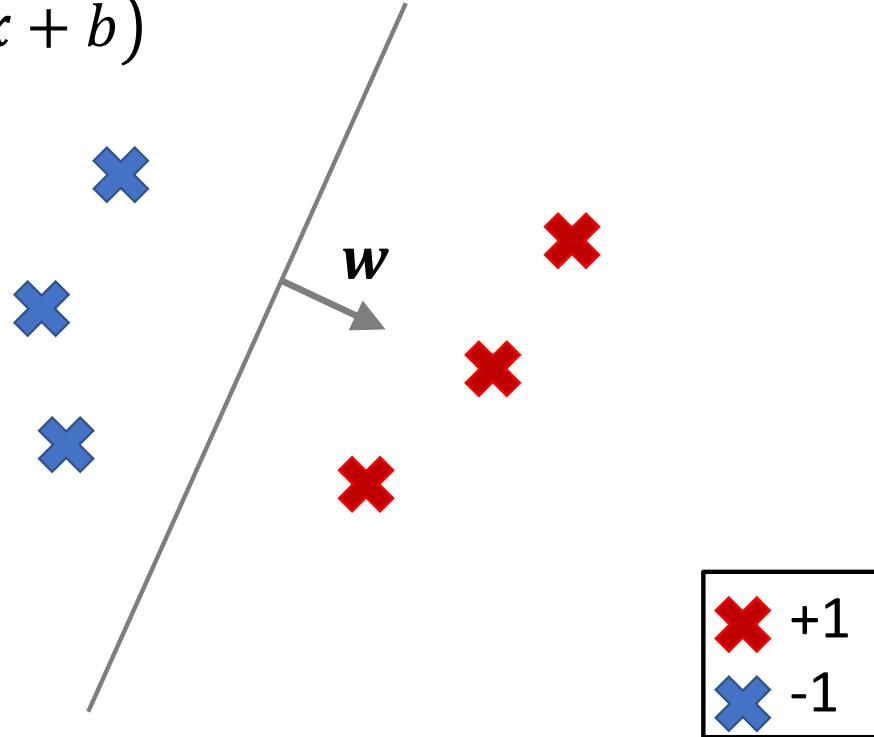
Linear classifier for binary classes

- To learn (\mathbf{w}, b)
- Decision rule: $\text{sign}(\mathbf{w}^T \mathbf{x} + b)$



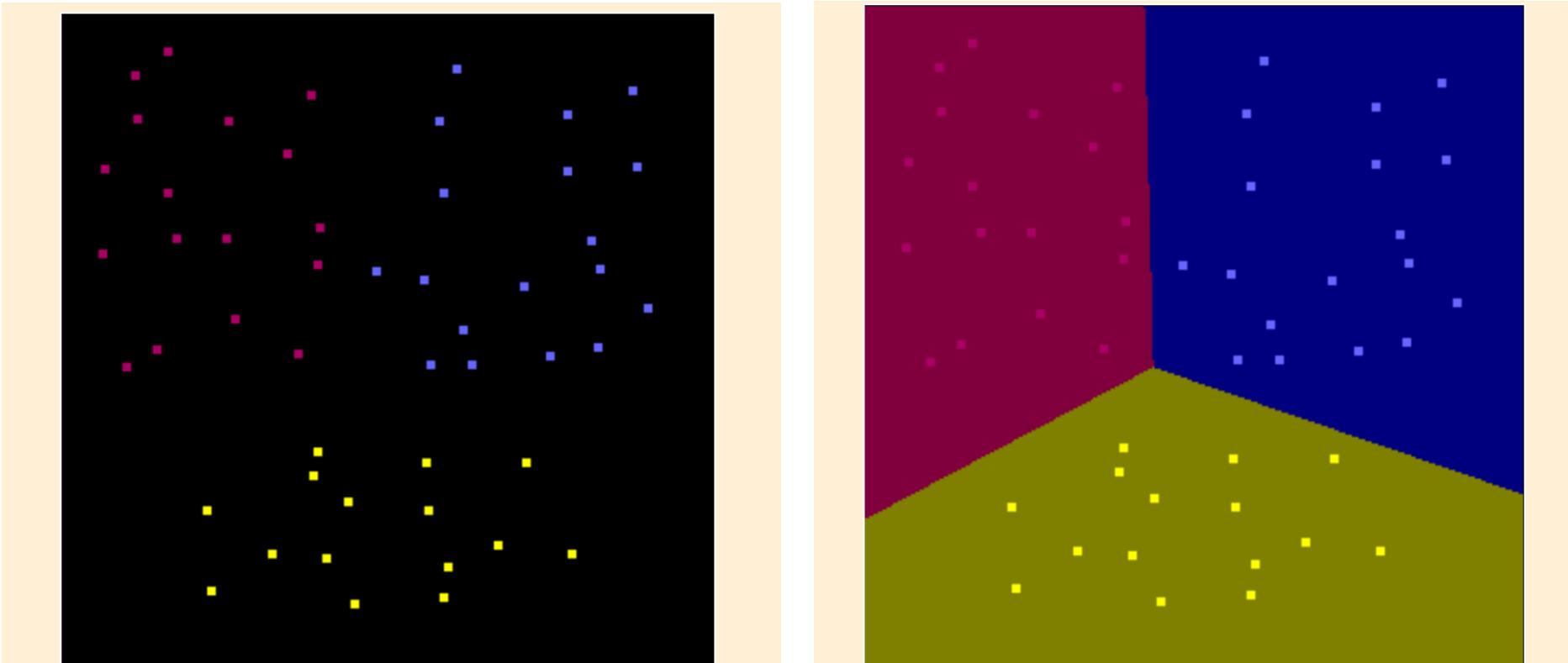
Linear classifier for binary classes

- To learn (w, b)
- Decision rule: $\text{sign}(w^T x + b)$



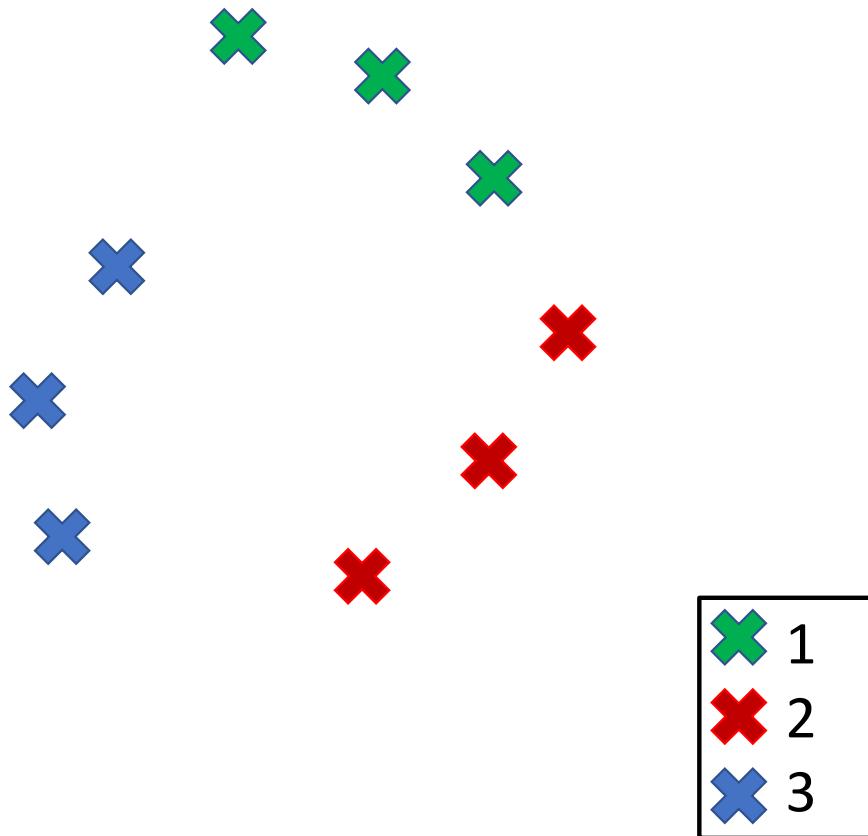
Classification

- To separate the “feature space” into class territories

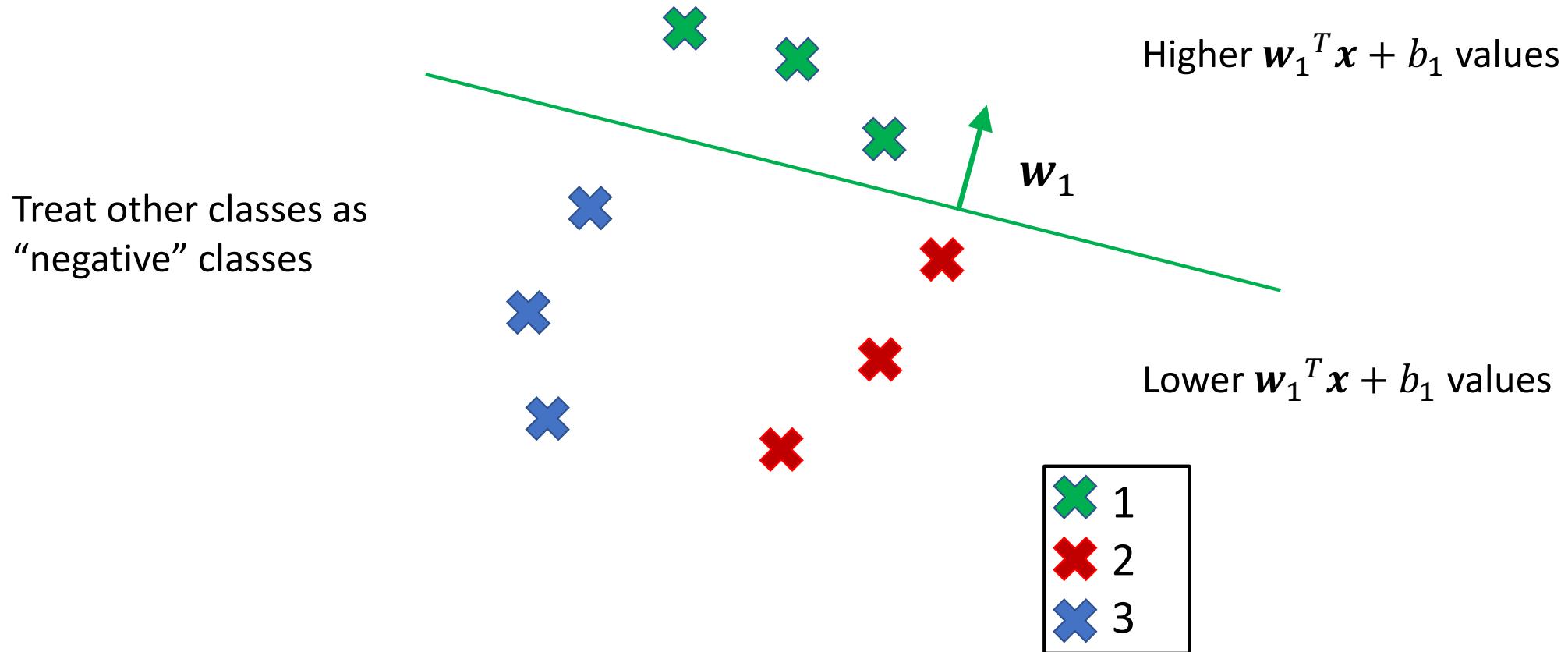


Linear classifier for multiple classes

- Linear separable?

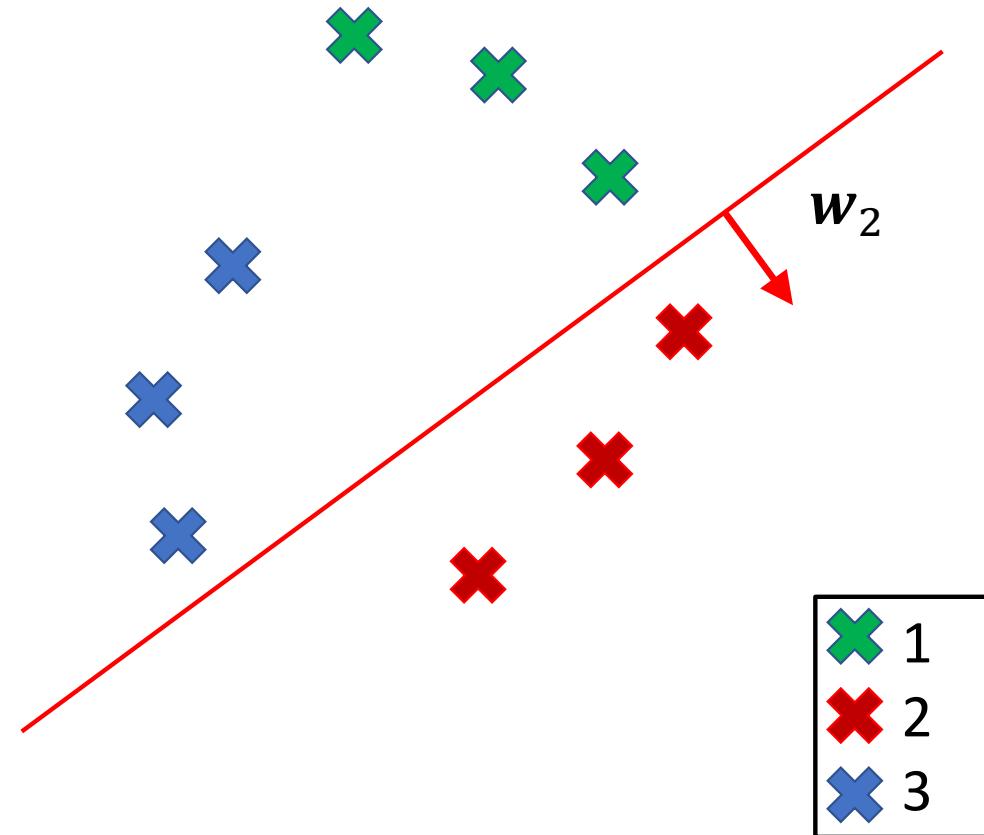


Linear classifier for multiple classes



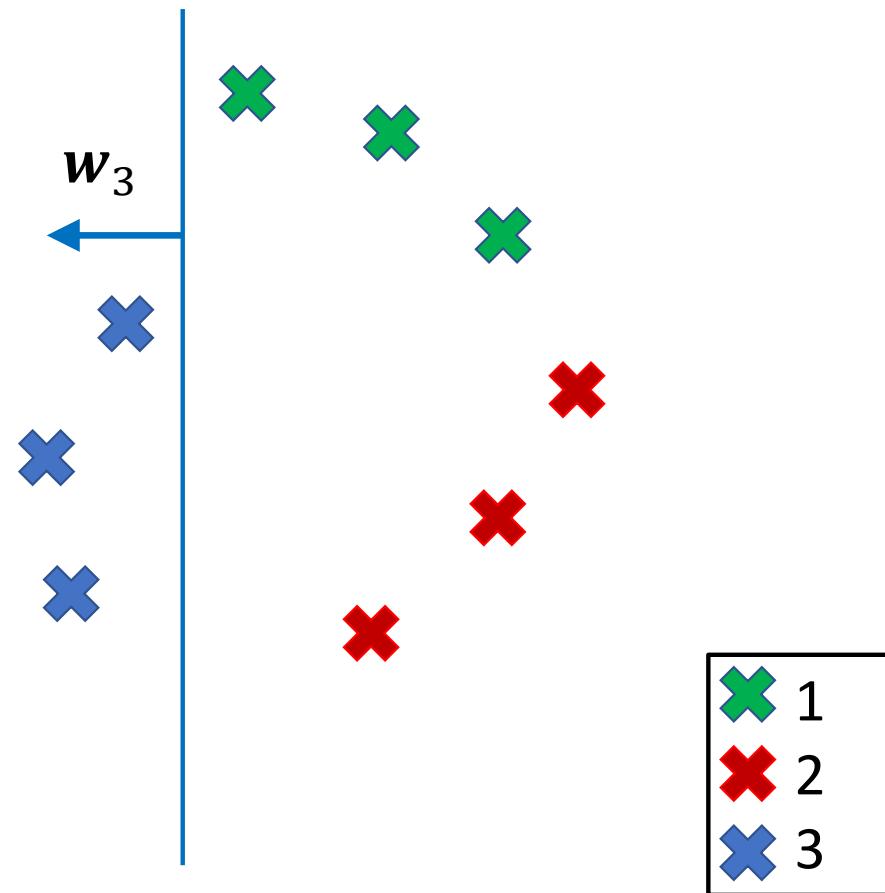
Linear classifier for multiple classes

Treat other classes as
“negative” classes



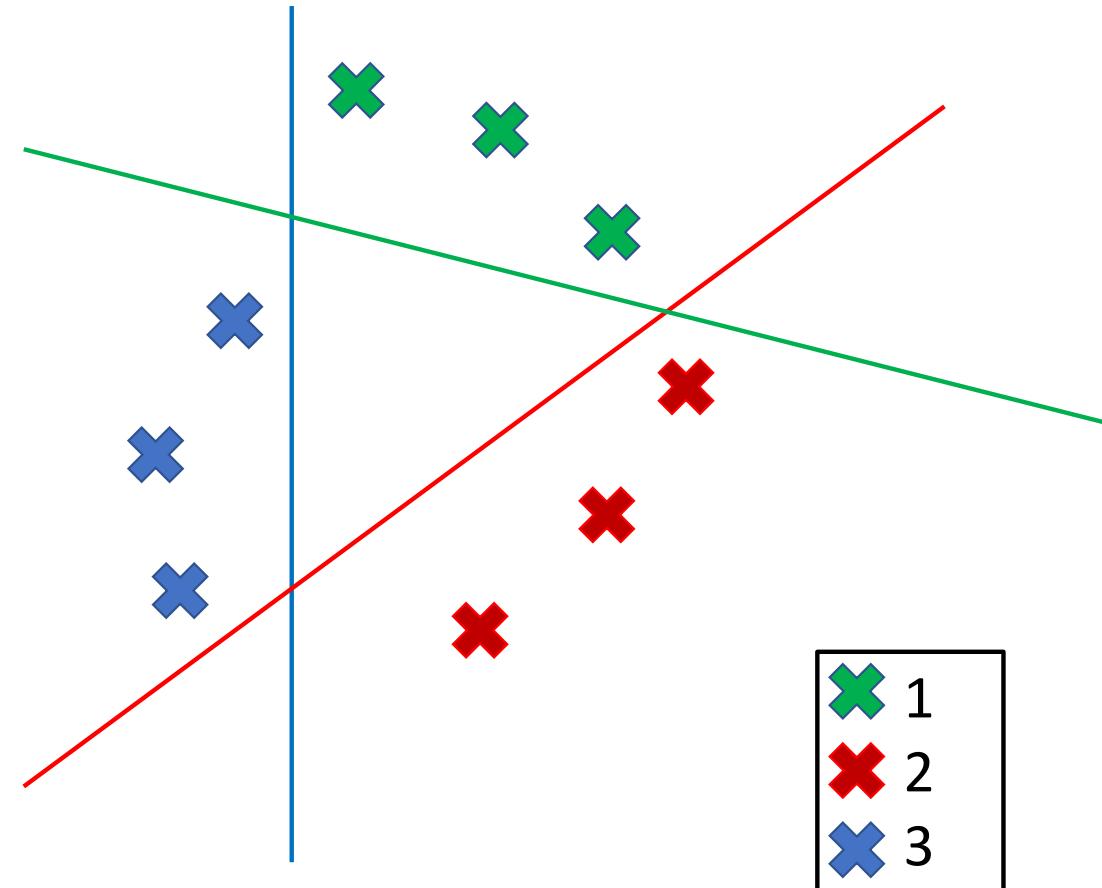
Linear classifier for multiple classes

Treat other classes as
“negative” classes



Linear classifier for multiple classes

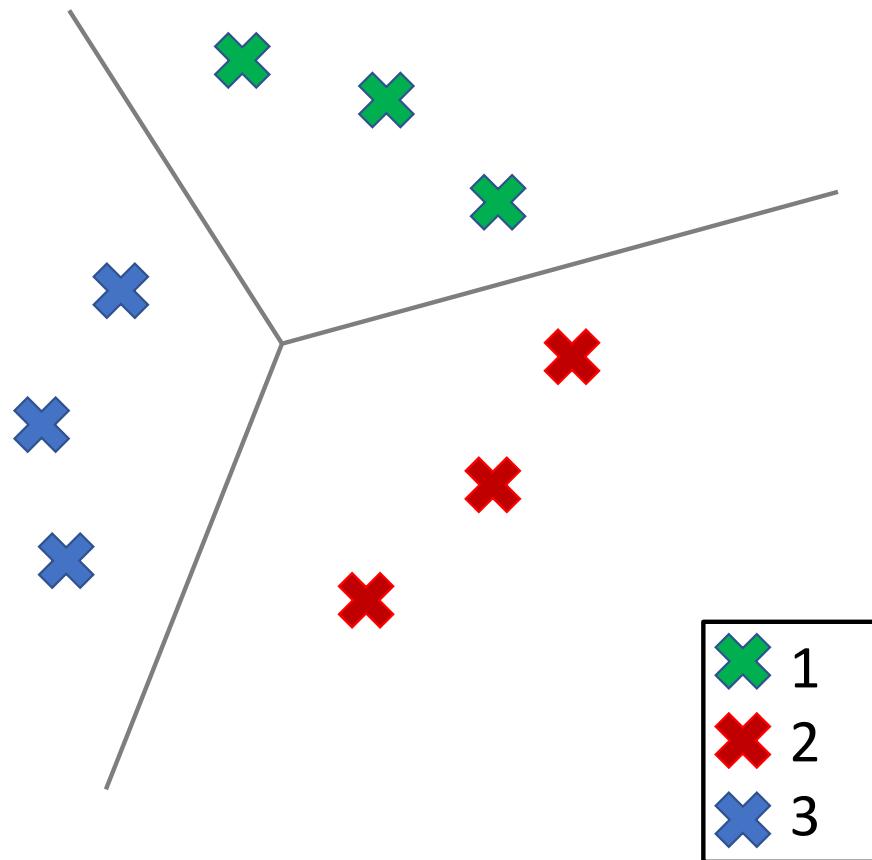
Every class c has its associated (w_c, b_c) that aims to give data from class c higher $w_c^T x + b_c$ values



Linear classifier for multiple classes

Decision rule:

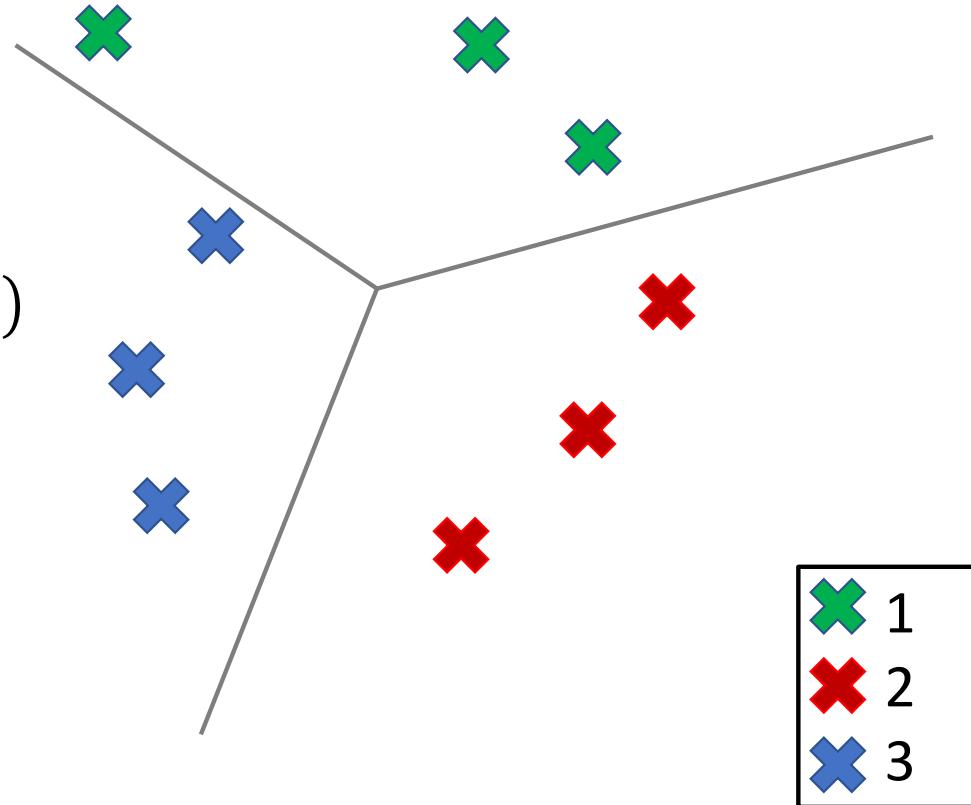
$$\hat{y} = \operatorname{argmax}_c (\mathbf{w}_c^T \mathbf{x} + b_c)$$



Linear classifier for multiple classes

Decision rule:

$$\hat{y} = \operatorname{argmax}_c (\mathbf{w}_c^T \mathbf{x} + b_c)$$

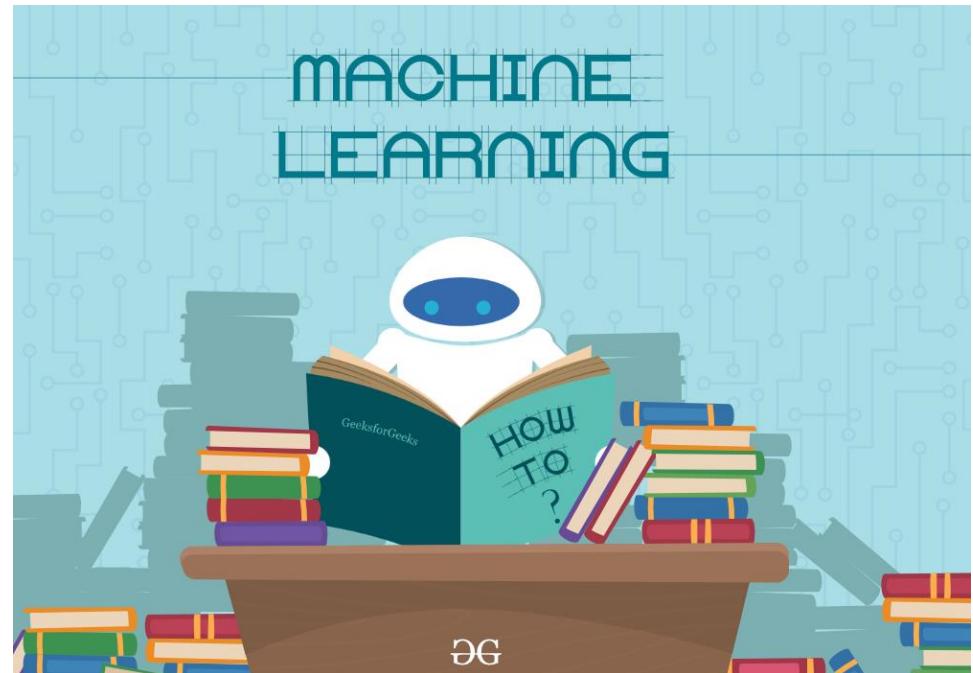


Today

Multi-class classification

- Basics
- Modeling + Learning
- Extension

ERM + regularization



Multi-class learning

- Training data: $D_{tr} = \{ (\mathbf{x}_n \in \mathbb{R}^D, y_n) \}_{i=1}^N$
 - $y \in \{1, \dots, C\}$. (The integer C here is not the hyperparameter in SVM.)
 - $y = \mathbf{e}_c = [0, \dots, 1, \dots, 0]^T$, a C -dim **one-hot** vector with only its c -th element = 1

Multi-class learning

- Training data: $D_{tr} = \{ (\mathbf{x}_n \in \mathbb{R}^D, y_n) \}_{i=1}^N$
 - $y \in \{1, \dots, C\}$. (The integer C here is not the hyperparameter in SVM.)
 - $y = \mathbf{e}_c = [0, \dots, 1, \dots, 0]^T$, a C -dim **one-hot** vector with only its c -th element = 1
- Softmax regression
- One-vs-all (or one-vs-other) formulation: extension of SVM, logistic regression
- Crammer-Singer SVM
- Other variants: one-vs-one, DAG (Directed Acyclic Graph)

Multi-class learning

- Training data: $D_{tr} = \{ (\mathbf{x}_n \in \mathbb{R}^D, y_n) \}_{i=1}^N$
 - $y \in \{1, \dots, C\}$. (The integer C here is not the hyperparameter in SVM.)
 - $y = \mathbf{e}_c = [0, \dots, 1, \dots, 0]^T$, a C -dim **one-hot** vector with only its c -th element = 1
- Softmax regression
- One-vs-all (or one-vs-other) formulation: extension of SVM, logistic regression
- Crammer-Singer SVM
- Other variants: one-vs-one, DAG (Directed Acyclic Graph)
- **Cautions: several “notations” (like λ, e_c) may be used in different contexts**

Softmax regression (cf. logistic regression)

- **Training data:** $D_{tr} = \{ (x_n \in \mathbb{R}^D, y_n \in \{1, \dots, C\}) \}_{i=1}^N$
- **Discriminative learning:** to construct $p(Y = c|x)$

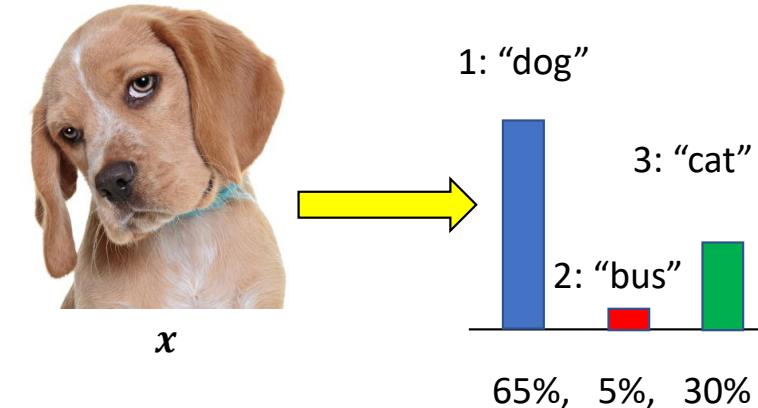
Softmax regression (cf. logistic regression)

- **Training data:** $D_{tr} = \{ (x_n \in \mathbb{R}^D, y_n \in \{1, \dots, C\}) \}_{i=1}^N$

- **Discriminative learning: to construct $p(Y = c|x)$**

- Conditional categorical; i.e., $p(Y = c|x) = \lambda(x)[c]$

- Modeling: $\lambda(x)[c] = \frac{e^{(\mathbf{w}_c^T x + b_c)}}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T x + b_{c'})}}$



Softmax regression (cf. logistic regression)

- **Training data:** $D_{tr} = \{ (x_n \in \mathbb{R}^D, y_n \in \{1, \dots, C\}) \}_{i=1}^N$

- **Discriminative learning: to construct $p(Y = c|x)$**

- Conditional categorical; i.e., $p(Y = c|x) = \lambda(x)[c]$

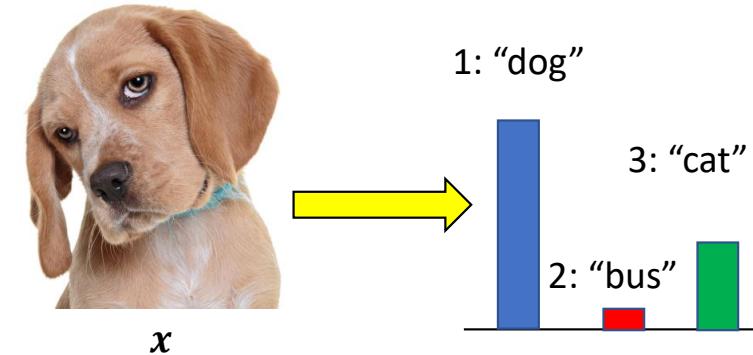
- Modeling: $\lambda(x)[c] = \frac{e^{(\mathbf{w}_c^T \mathbf{x} + b_c)}}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x} + b_{c'})}}$

- Properties 1: $1 \geq \lambda(x)[c] \geq 0, \sum_{c=1}^C \lambda(x)[c] = 1$

- Properties 2: $\mathbf{w}_c^T \mathbf{x} + b_c \uparrow, p(Y = c|x) \uparrow$

- Properties 3: $\mathbf{w}_c^T \mathbf{x} + b_c$ is usually called the decision value or logit

- Question: what if $C = 2$? We get \mathbf{w}_1 and \mathbf{w}_2 but logistic regression only has \mathbf{w} (HW # 4)



Softmax regression

- **Learning (MLE):** $\widehat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(D_{tr}; \boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \prod_n p(y_n | x_n; \boldsymbol{\theta}) p(x_n)$

Softmax regression

- **Learning (MLE):** $\widehat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(D_{tr}; \boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \prod_n p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) p(\mathbf{x}_n)$

$$\begin{aligned} \log \prod_{n=1}^N p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) &= \sum_{n=1}^N \log p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) = \sum_{n=1}^N \log \frac{e^{(\mathbf{w}_{y_n}^T \mathbf{x}_n)}}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)}} , \text{ ignore } b_c \\ &= \sum_{n=1}^N \left[\mathbf{w}_{y_n}^T \mathbf{x}_n - \log \sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)} \right] \end{aligned}$$

Softmax regression

- **Learning (MLE):** $\widehat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(D_{tr}; \boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_n p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) p(\mathbf{x}_n)$

$$\begin{aligned} \log \prod_{n=1}^N p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) &= \sum_{n=1}^N \log p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) = \sum_{n=1}^N \log \frac{e^{(\mathbf{w}_{y_n}^T \mathbf{x}_n)}}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)}} , \text{ ignore } b_c \\ &= \sum_{n=1}^N \left[\mathbf{w}_{y_n}^T \mathbf{x}_n - \log \sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)} \right] \end{aligned}$$

- $\{\mathbf{w}_c\}_{c=1}^C \stackrel{\text{MLE}}{=} \underset{\{\mathbf{w}_c\}_{c=1}^C}{\operatorname{argmin}} \sum_{n=1}^N \left[-\mathbf{w}_{y_n}^T \mathbf{x}_n + \log \sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)} \right]$

Softmax regression: gradient

- Gradient for each data instance (x_n, y_n)

$$\nabla_{\mathbf{w}_c} \left[-\mathbf{w}_{y_n}^T \mathbf{x}_n + \log \sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)} \right]$$

Softmax regression: gradient

- Gradient for each data instance (\mathbf{x}_n, y_n)

$$\begin{aligned}\nabla_{\mathbf{w}_{\textcolor{red}{c}}}\left[-\mathbf{w}_{y_n}^T \mathbf{x}_n + \log \sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)}\right] &= -\mathbf{1}[y_n == c] \times \mathbf{x}_n + \frac{\nabla_{\mathbf{w}_{\textcolor{red}{c}}}\left[\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)}\right]}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)}} \\ &= -\mathbf{1}[y_n == c] \times \mathbf{x}_n + \frac{\nabla_{\mathbf{w}_{\textcolor{red}{c}}} e^{(\mathbf{w}_{\textcolor{red}{c}}^T \mathbf{x}_n)}}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)}} = -\mathbf{1}[y_n == c] \times \mathbf{x}_n + \frac{e^{(\mathbf{w}_{\textcolor{red}{c}}^T \mathbf{x}_n)} \mathbf{x}_n}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x}_n)}} \\ &= -\mathbf{1}[y_n == c] \times \mathbf{x}_n + p(c|\mathbf{x}_n; \mathbf{w}) \mathbf{x}_n = -\{\mathbf{1}[y_n == c] - p(c|\mathbf{x}_n; \mathbf{w})\} \mathbf{x}_n\end{aligned}$$

- Note: Every (\mathbf{x}_n, y_n) contributes to the update of all $\mathbf{w}_c, \forall c \in \{1, \dots, C\}$, not only to \mathbf{w}_{y_n}

Softmax regression: cross entropy loss

- **Learning (MLE):** $\widehat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(D_{tr}; \boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \prod_n p(y_n | x_n; \boldsymbol{\theta}) p(x_n)$

$$-\log \prod_{n=1}^N p(y_n | x_n; \boldsymbol{\theta}) = -\sum_{n=1}^N \log p(y_n | x_n; \boldsymbol{\theta})$$

Softmax regression: cross entropy loss

- **Learning (MLE):** $\widehat{\boldsymbol{\theta}} = \operatorname{argmax}_{\boldsymbol{\theta}} p(D_{tr}; \boldsymbol{\theta}) = \operatorname{argmax}_{\boldsymbol{\theta}} \prod_n p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) p(\mathbf{x}_n)$

$$\begin{aligned} -\log \prod_{n=1}^N p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) &= -\sum_{n=1}^N \log p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) \\ &= \sum_{n=1}^N \sum_{c=1}^C -\mathbf{1}[y_n == c] \times \log p(c | \mathbf{x}_n; \boldsymbol{\theta}) \end{aligned}$$

Softmax regression: cross entropy loss

- **Learning (MLE):** $\widehat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} p(D_{tr}; \boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_n p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) p(\mathbf{x}_n)$

$$\begin{aligned} -\log \prod_{n=1}^N p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) &= -\sum_{n=1}^N \log p(y_n | \mathbf{x}_n; \boldsymbol{\theta}) \\ &= \sum_{n=1}^N \sum_{c=1}^C -\mathbf{1}[y_n == c] \times \log p(c | \mathbf{x}_n; \boldsymbol{\theta}) \\ &= \sum_{n=1}^N \sum_{c=1}^C -\mathbf{e}_{y_n}[c] \times \log p(c | \mathbf{x}_n; \boldsymbol{\theta}) = -\sum_{n=1}^N \mathbf{e}_{y_n}^T \begin{bmatrix} \log p(1 | \mathbf{x}_n; \boldsymbol{\theta}) \\ \vdots \\ \log p(C | \mathbf{x}_n; \boldsymbol{\theta}) \end{bmatrix} \end{aligned}$$

Softmax regression: cross entropy loss

- Learning (MLE): $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} p(D_{tr}; \theta) = \underset{\theta}{\operatorname{argmax}} \prod_n p(y_n | x_n; \theta) p(x_n)$

$$-\log \prod_{n=1}^N p(y_n | x_n; \theta) = -\sum_{n=1}^N \log p(y_n | x_n; \theta)$$

$$= \sum_{n=1}^N \sum_{c=1}^C -\mathbf{1}[y_n == c] \times \log p(c | x_n; \theta)$$

$$= \sum_{n=1}^N \sum_{c=1}^C -\boxed{e_{y_n}[c]} \times \boxed{\log p(c | x_n; \theta)} = -\sum_{n=1}^N e_{y_n}^T \begin{bmatrix} \log p(1 | x_n; \theta) \\ \vdots \\ \log p(C | x_n; \theta) \end{bmatrix}$$

target prediction

Softmax regression: cross entropy loss

- Learning (MLE): $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} p(D_{tr}; \theta) = \underset{\theta}{\operatorname{argmax}} \prod_n p(y_n | x_n; \theta) p(x_n)$

$$\begin{aligned}& -\log \prod_{n=1}^N p(y_n | x_n; \theta) = -\sum_{n=1}^N \log p(y_n | x_n; \theta) \\&= \sum_{n=1}^N \sum_{c=1}^C -\mathbf{1}[y_n == c] \times \log p(c | x_n; \theta) \\&= \sum_{n=1}^N \sum_{c=1}^C -\boxed{\mathbf{e}_{y_n}[c]} \times \boxed{\log p(c | x_n; \theta)} = -\sum_{n=1}^N \mathbf{e}_{y_n}^T \begin{bmatrix} \log p(1 | x_n; \theta) \\ \vdots \\ \log p(C | x_n; \theta) \end{bmatrix}\end{aligned}$$

target prediction

- Re-write the objective
- Want the target and the prediction to be similar (i.e., to have high inner product)

Softmax regression: prediction

- Prediction by Bayes optimal classifier

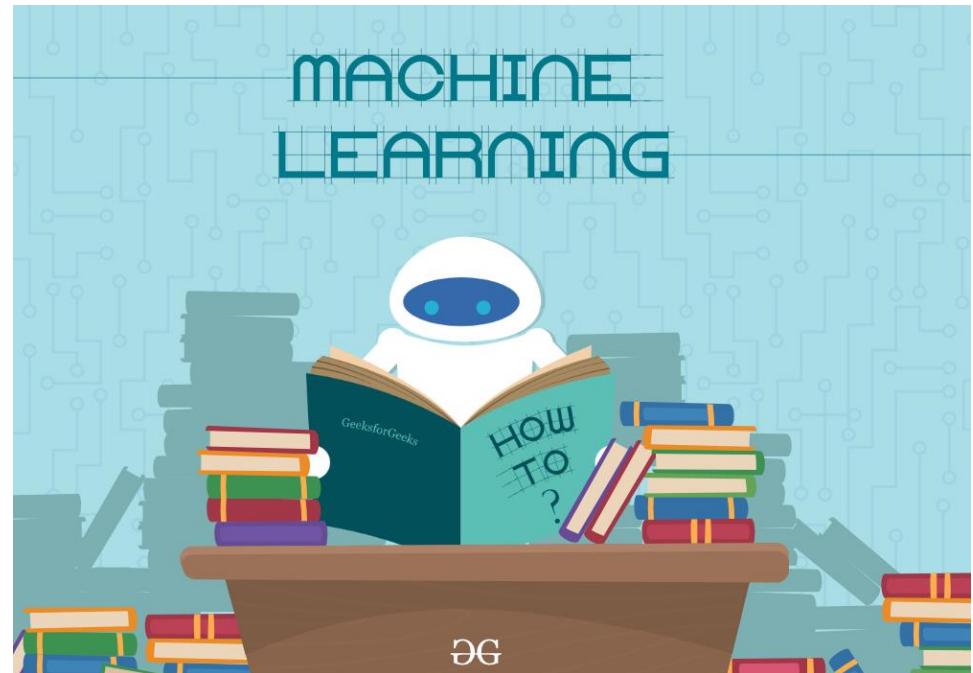
$$\circ \hat{y} = \operatorname{argmax}_c p(c|x; \theta) = \operatorname{argmax}_c \frac{e^{(\mathbf{w}_{\textcolor{red}{c}}^T \mathbf{x})}}{\sum_{c'=1}^C e^{(\mathbf{w}_{c'}^T \mathbf{x})}} = \operatorname{argmax}_c e^{(\mathbf{w}_{\textcolor{red}{c}}^T \mathbf{x})} = \operatorname{argmax}_c \mathbf{w}_{\textcolor{red}{c}}^T \mathbf{x}$$

Today

Multi-class classification

- Basics
- Modeling + Learning
- Extension

ERM + regularization



One-vs-all (or one-vs-other)



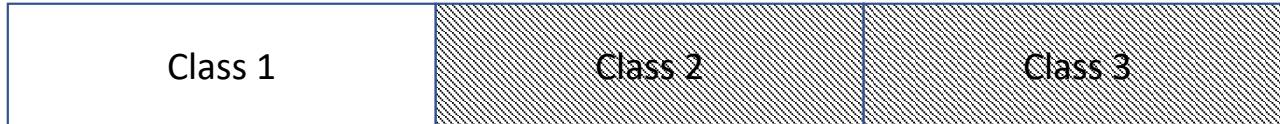
- One-vs-all learns w_c by viewing data of class c as class +1, others as class -1

One-vs-all (or one-vs-other)



- One-vs-all learns w_c by viewing data of class c as class +1, others as class -1

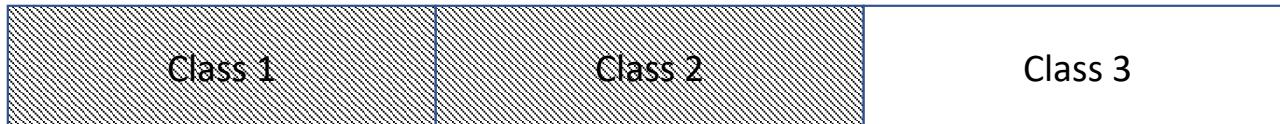
- w_1 :



- w_2 :



- w_3 :



One-vs-all (or one-vs-other)

- Learn each \mathbf{w}_c independently
- For $c \in \{1, \dots, C\}$
 - Temporally set $y_n = \begin{cases} +1, & \text{if } y_n = c \\ -1, & \text{otherwise} \end{cases}$, such that $D_{tr} = \{(\mathbf{x}_n \in \mathbb{R}^D, y_n \in \{+1, -1\})\}_{i=1}^N$
 - \mathbf{w}_c is the learned \mathbf{w} from binary classification (e.g., by SVM, logistic regression)

One-vs-all (or one-vs-other)

- Learn each \mathbf{w}_c independently
- For $c \in \{1, \dots, C\}$
 - Temporally set $y_n = \begin{cases} +1, & \text{if } y_n = c \\ -1, & \text{otherwise} \end{cases}$, such that $D_{tr} = \{(\mathbf{x}_n \in \mathbb{R}^D, y_n \in \{+1, -1\})\}_{i=1}^N$
 - \mathbf{w}_c is the learned \mathbf{w} from binary classification (e.g., by SVM, logistic regression)
- Prediction: $\operatorname{argmax}_c \mathbf{w}_c^T \mathbf{x}$
 - In training, we encourage $\mathbf{w}_c^T \mathbf{x}_n$ to be large if $y_n = c$; small if $y_n \neq c$
 - Problem: no calibration among \mathbf{w}_c
- Note: multi-class logistic regression vs. softmax regression

Crammer-singer SVM (learn w_c together)

- **Prediction:** $\operatorname{argmax}_c w_c^T x$
 - That is, if x belongs to class y , we hope that $w_y^T x > \max_{c \neq y} w_c^T x$
- **Learning:**
 - $\min \frac{1}{2} \sum_c w_c^T w_c + \lambda \sum_n \max \left\{ 0, \left(\max_{c \neq y_n} w_c^T x_n + 1 \right) - w_{y_n}^T x_n \right\}$
 - $\min \frac{1}{2} \sum_c w_c^T w_c + \lambda \sum_n \max \left\{ 0, \left(\max_c w_c^T x_n + \mathbf{1}[y_n \neq c] \right) - w_{y_n}^T x_n \right\}$
 - The 1 or $\mathbf{1}[y_n \neq c]$ here is the enforced margin
 - Here I use λ to denote the hyper-parameter

Questions?



THE OHIO STATE UNIVERSITY

Other variants

Class 1

Class 2

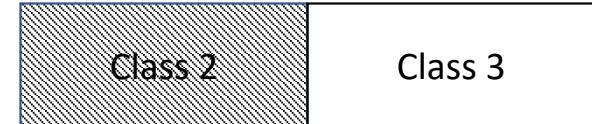
Class 3

- **One-vs-one learning:**

- Build $C \times (C - 1)/2$ binary classifiers for pairwise comparison
- Ex: w_{1-2} :



w_{3-2} :



- **Prediction:**

- Accumulate the decision values of each class, for predicting the label for x
 - Ex: $\text{score}(c) = \sum_{c' \neq c} (w_{c-c'}^T x)$
- Comparison along a DAG (Directed Acyclic Graph)

Comparison

- Modeling and prediction:
 - $\underset{c}{\operatorname{argmax}} \mathbf{w}_c^T \mathbf{x}$: softmax regression, Crammer-Singer, one-vs-all
 - One-vs-one $\mathbf{w}_{c-c'}$: $\underset{c}{\operatorname{argmax}} \sum_{c' \neq c} (\mathbf{w}_{c-c'}^T \mathbf{x})$ or DAG or other accumulation rule
- Learning:
 - Learn $\{\mathbf{w}_c\}_{c=1}^C$ jointly: softmax regression, Crammer-Singer
 - Learn $\{\mathbf{w}_c\}_{c=1}^C$ or $\{\mathbf{w}_{c-c'}\}_{c=1, c'=1}^{C,C}$ independently: one-vs-all, one-vs-one
- Problems for one-vs-all and one-vs-one:
 - Limited or imbalanced data (too many -1 data than +1 data)

Summary

- Modeling
 - **Decision values:** Linear, nonlinear, kernelized, nearest neighbor, ...
 - $w_c^T x + b_c$
 - $w_c^T \phi(x)$
 - $\sum_n \alpha_{c,n} k_\phi(x_n, x)$
 - **Parametric vs. nonparametric:** Does the number of parameters grow with training data?
- Learning
 - Empirical risk minimization (ERM)
 - Structured risk minimization, regularized risk/loss minimization
 - Probabilistic: MLE, MAP

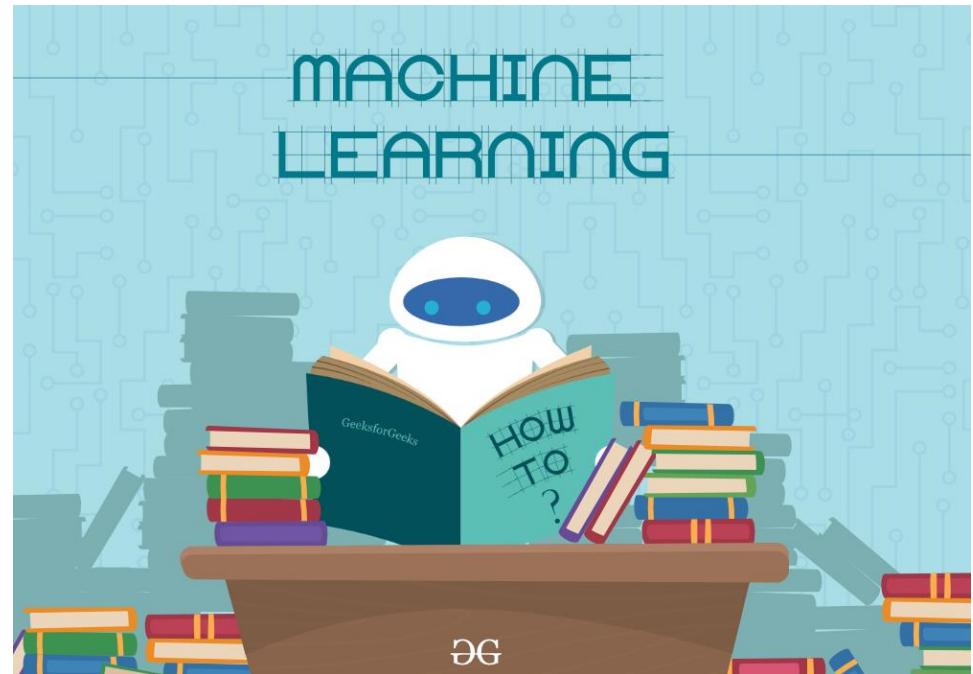
Today

Multi-class classification

- Basics
- Modeling + Learning
- Extension

ERM + regularization

- Warm up
- Basic



Logistic regression: another perspective

- **Training data:** $D_{tr} = \{ (x_i \in \mathbb{R}^D, y_i \in \{+1, -1\}) \}_{i=1}^N$
- Learning: $\tilde{\mathbf{w}}^{\text{MLE}} = \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_{i=1}^N \log \left(1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{x}_i} \right)$
- Prediction: $\operatorname{sign}(\mathbf{w}^T \mathbf{x} + b)$

Logistic regression: another perspective

- Training data: $D_{tr} = \{ (x_i \in \mathbb{R}^D, y_i \in \{+1, -1\}) \}_{i=1}^N$
 - Learning: $\tilde{\mathbf{w}}^{\text{MLE}} = \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_{i=1}^N \log \left(1 + e^{-y_i \tilde{\mathbf{w}}^T \tilde{x}_i} \right) =$
$$\underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_{i=1}^N \ell(\tilde{\mathbf{w}}^T \tilde{x}_i, y_i) = \underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_{i=1}^N \ell(h_{\tilde{\mathbf{w}}}(\tilde{x}_i), y_i) =$$
$$\underset{\tilde{\mathbf{w}}}{\operatorname{argmin}} \sum_{i=1}^N \ell(\hat{y}_i, y_i)$$
 - Prediction: $\operatorname{sign}(\mathbf{w}^T \mathbf{x} + b)$
- “soft” prediction! Loss function

Empirical risk minimization recap

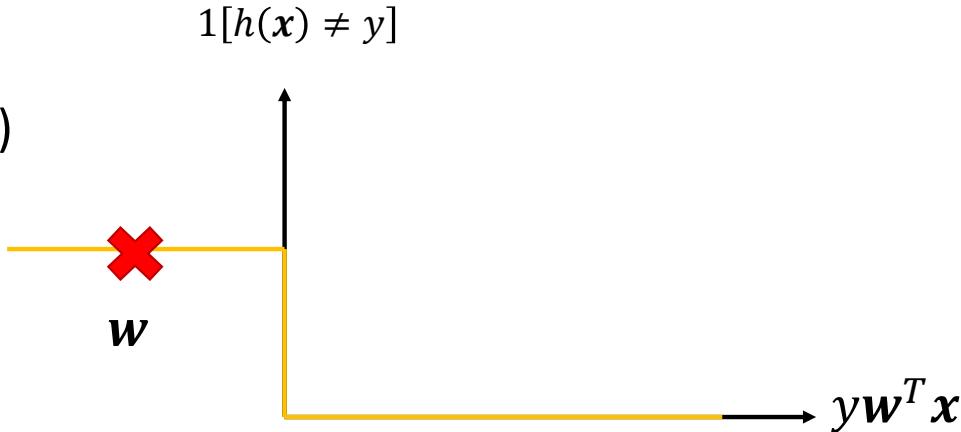
- In order to achieve a low generalization loss $E_{(x,y) \sim P}[\ell(h(x), y)]$
 - We minimize the training loss + certain regularization
 - Training loss: $\frac{1}{N} \sum_{n=1}^N \ell(h(x_n), y_n)$, where $D_{tr} = \{(x_n \in \mathbb{R}^D, y_n \in \{+1, -1\})\}_{i=1}^N$
 - Ideally, h & ℓ in the training should be the same as h & ℓ in the test (generalization) loss

Empirical risk minimization recap

- In order to achieve a low generalization loss $E_{(x,y) \sim P}[\ell(h(x), y)]$
 - We minimize the training loss + certain regularization
 - Training loss: $\frac{1}{N} \sum_{n=1}^N \ell(h(x_n), y_n)$, where $D_{tr} = \{(x_n \in \mathbb{R}^D, y_n \in \{+1, -1\})\}_{i=1}^N$
 - Ideally, h & ℓ in the training should be the same as h & ℓ in the test (generalization) loss
- For classification, we care
 - 0-1 loss: $\ell(y, \hat{y}) = 1[y \neq \hat{y}]$; $1[\text{True}] = 1, 1[\text{False}] = 0$ is called an indicator function
- For linear classifier, we use
 - $\hat{y} = h(x) = \text{sign}(w^T x)$; the bias term b is absorbed into w for brevity

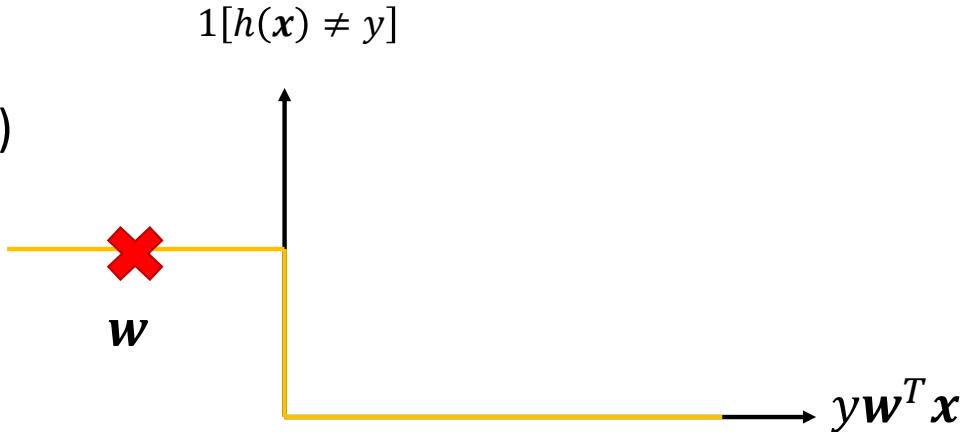
Problems with 0-1 losses

- $\ell(y, \hat{y}) = 1[y \neq \hat{y}]$
 - Not differentiable w.r.t. \hat{y} (i.e., the prediction)
 - Example: for $y_i \in \{+1, -1\}$, we want
 - \hat{y} (and thus $w^T x$) and y to be of the same sign



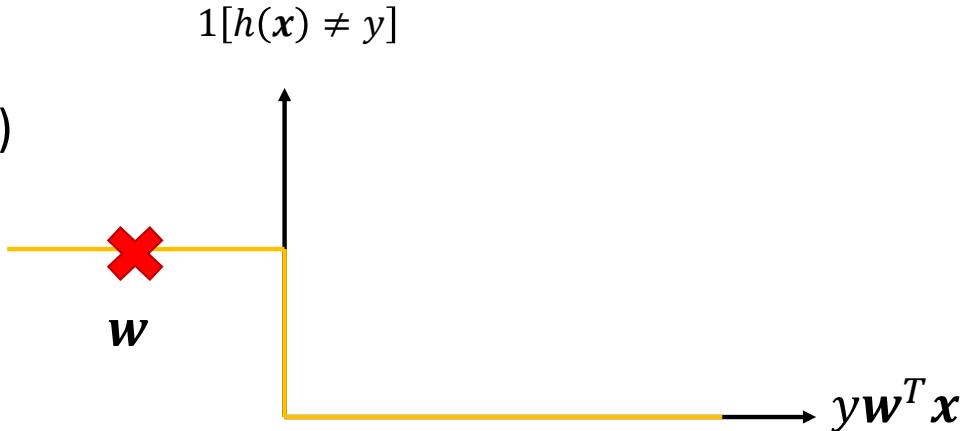
Problems with 0-1 losses

- $\ell(y, \hat{y}) = 1[y \neq \hat{y}]$
 - Not differentiable w.r.t. \hat{y} (i.e., the prediction)
 - Example: for $y_i \in \{+1, -1\}$, we want
 - \hat{y} (and thus $w^T x$) and y to be of the same sign
 - Hard to optimize w
 - No signals of where w should move



Problems with 0-1 losses

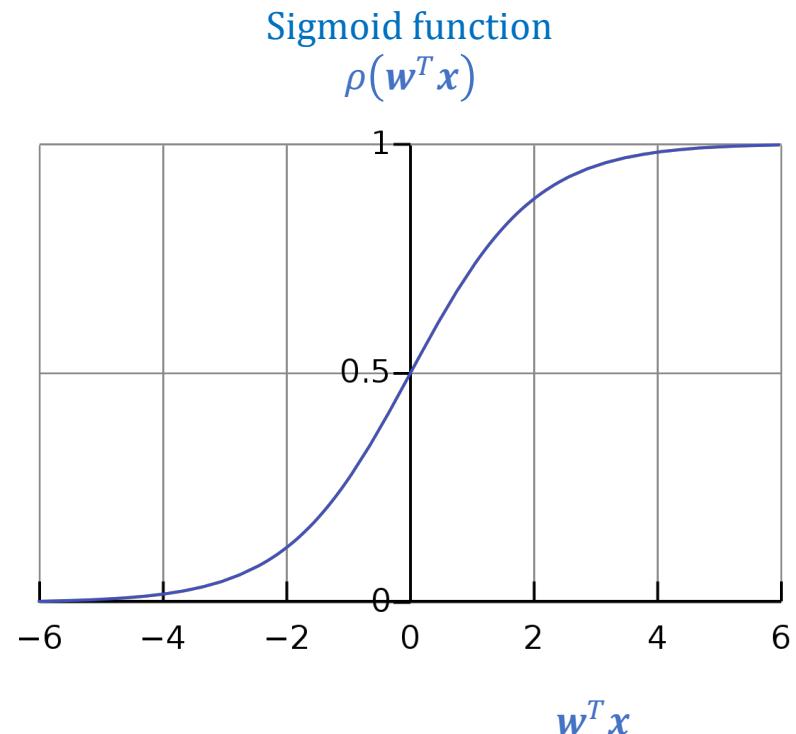
- $\ell(y, \hat{y}) = 1[y \neq \hat{y}]$
 - Not differentiable w.r.t. \hat{y} (i.e., the prediction)
 - Example: for $y_i \in \{+1, -1\}$, we want
 - \hat{y} (and thus $w^T x$) and y to be of the same sign
 - Hard to optimize w
 - No signals of where w should move



- We need a relaxed (surrogate) loss “**in training**” to encourage:
 - $y = 1, \quad w^T x \uparrow$
 - $y = -1, \quad w^T x \downarrow$

Logistic loss

- Let $p(+1|x; \mathbf{w}) = \rho(\mathbf{w}^T \mathbf{x})$
- Let $p(-1|x; \mathbf{w}) = 1 - \rho(\mathbf{w}^T \mathbf{x})$
- Want to maximize: $p(y_n|x_n; \mathbf{w})$
- Want to minimize: $-\log p(y_n|x_n; \mathbf{w})$
 $= -\log \rho(y_n \mathbf{w}^T \mathbf{x})$
- When
 - $y = 1, \mathbf{w}^T \mathbf{x} \uparrow$
 - $y = -1, \mathbf{w}^T \mathbf{x} \downarrow$

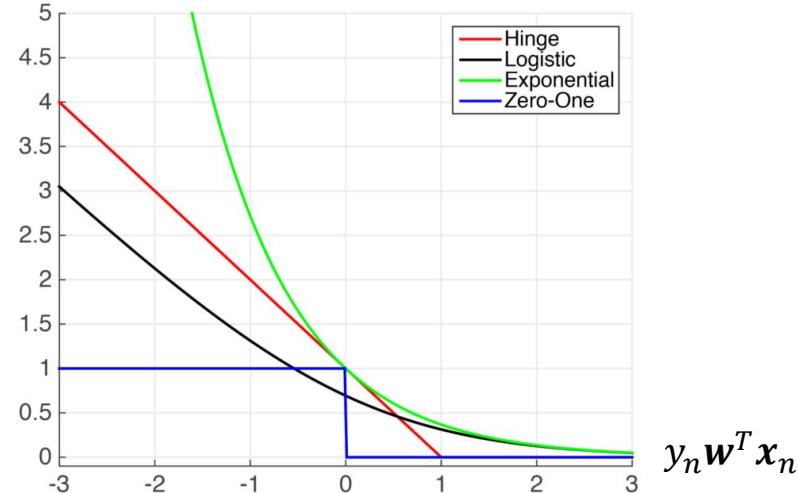


Questions!

- What properties must a surrogate loss have?
- Are there any other loss functions?
- Are there any other regularization functions, beyond $\lambda\|w\|_2^2$?

A first glance

- What properties must a surrogate loss have?
 - Upper bounding the 0-1 loss: so, small surrogate loss implies small 0-1 loss
- Are there any other loss functions?



- Are there any other regularization functions, beyond $\lambda \|\mathbf{w}\|_2^2$?
 - Yes: as long as it is “lower bounded” and characterizes model complexity

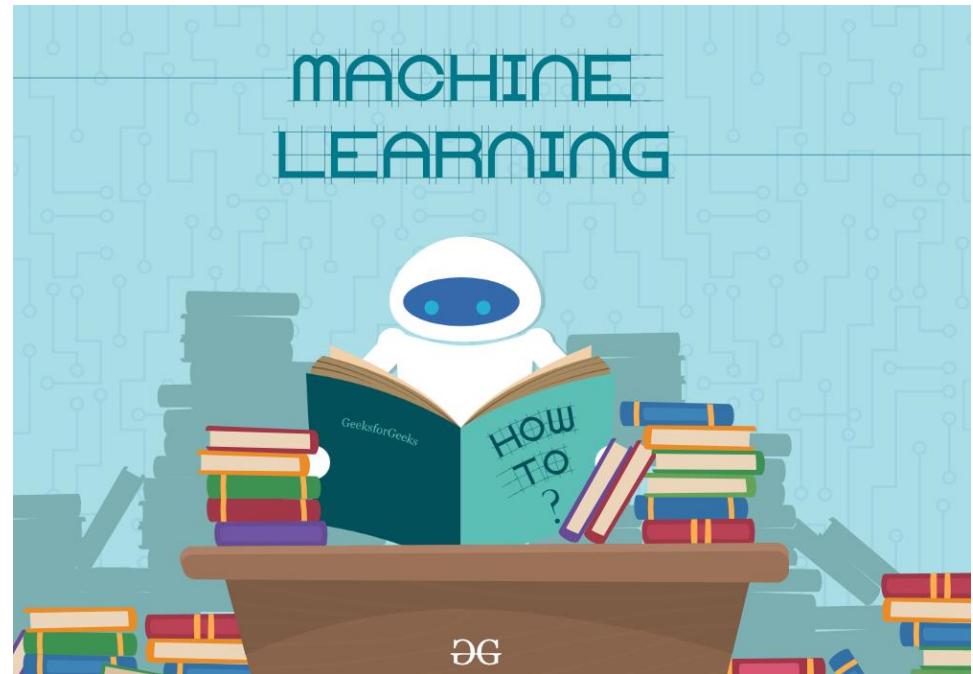
Today

Multi-class classification

- Basics
- Modeling + Learning
- Extension

ERM + regularization

- Warm up
- Basic



Empirical risk minimization

- **Goal:** find $h: \mathcal{X} \mapsto \mathcal{Y}$ to minimize $E_{(x,y) \sim P}[\ell_{te}(h(x), y)]$, where P is unknown
- What do we know about P ?
 - $D_{tr} = \{(x_n \in \mathbb{R}^D, y_n)\}_{i=1}^N$ that are sampled IID from P

Empirical risk minimization

- **Goal:** find $h: \mathcal{X} \mapsto \mathcal{Y}$ to minimize $E_{(x,y) \sim P}[\ell_{te}(h(x), y)]$, where P is unknown
- What do we know about P ?

○ $D_{tr} = \{(x_n \in \mathbb{R}^D, y_n)\}_{i=1}^N$ that are sampled IID from P

- **Empirical risk minimization (regularized risk/loss minimization)**

$$\hat{h} = \operatorname{argmin}_{h \in \mathcal{H}} \frac{1}{N} \sum_{n=1}^N \ell_{tr}(h(x_n), y_n) + \lambda \times r(h)$$

Training loss Regularization

- ℓ_{tr} : (continuous) function that penalizes training error
- r : (continuous) function that penalizes classifier/hypothesis complexity
- λ : hyper-parameter

Empirical risk minimization: example

- Linear soft SVM:

- $y \in \{+1, -1\}, h_{\mathbf{w}}(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$

- $\underset{\mathbf{w}}{\operatorname{argmin}} \left[\frac{1}{N} \sum_n [1 - y_n(\mathbf{w}^T \mathbf{x}_n)]_+ + \lambda \|\mathbf{w}\|_2^2 \right]$
Hinge-loss L2-regularization

- For homework implementation, please look at “**HW3.pdf**”

Empirical risk minimization: example

- Linear soft SVM:
 - $y \in \{+1, -1\}$, $h_{\mathbf{w}}(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x}_n)$
 - $\underset{\mathbf{w}}{\operatorname{argmin}} \quad \boxed{\frac{1}{N} \sum_n [1 - y_n (\mathbf{w}^T \mathbf{x}_n)]_+} + \boxed{\lambda \|\mathbf{w}\|_2^2}$

Hinge-lossL2-regularization
 - For homework implementation, please look at “**HW3.pdf**”
- Any other kinds of (surrogate) loss functions and regularization functions?
 - If so, what are the principles or requirements?