

10-701 Midterm Review

Mar.19, 2018

Model-based Approach

Data: X_1, X_2, \dots, X_n .

Model: $P(X; \theta)$ with parameters θ .

Assumption: Data drawn *i.i.d* from distribution $P(X; \theta^*)$ for some unknown θ^* .

Mission (should you choose to accept it): recover θ^* from data X_1, X_2, \dots, X_n .

- **Estimate Model Parameter**
 - Different ways: MLE/MAP
- Use model in production

Maximum Likelihood Estimation

Choose θ that maximizes the probability of observed data

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(D | \theta)$$

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(X_1, X_2, \dots, X_n; \theta)$$

$$= \arg \max_{\theta} \prod_{i=1}^n P(X_i; \theta)$$

$$= \arg \max_{\theta} \sum_{i=1}^n \log(P(X_i; \theta))$$

Take first derivative
Put to zero

How good is MLE?

- Asymptotically, unbiased
- Consistent: under some constraints
- Finite Sample:
 - Sample Question: flipped a coin once and observed a head.
 - What is MLE for p , where p is probability of heads showing.
 - Using MLE, give an interval such that there's a 95% chance of p lying in it.

MAP

Choose θ whose probability given data is highest,
ie maximize posterior distribution

$$\hat{\theta}_{MAP} = \operatorname{argmax}_{\theta} P(\theta | X_1, X_2, \dots, X_n)$$

$$P(\theta | \mathcal{D}) = \frac{P(\mathcal{D} | \theta)P(\theta)}{P(\mathcal{D})}$$

$$P(\theta | \mathcal{D}) \propto P(\mathcal{D} | \theta)P(\theta)$$

posterior likelihood prior

$$\begin{aligned}\hat{\theta}_{MAP} &= \operatorname{argmax}_{\theta} \left(\prod_{i=1}^n P(X_i; \theta) \right) P(\theta) \\ &= \operatorname{argmax}_{\theta} \sum_{i=1}^n \log(P(X_i; \theta)) + \log(P(\theta))\end{aligned}$$

Regularizer

Questions?

- Practice calculating MLE's and MAPs
 - Coin flips (prior: Beta)
 - Dice Flips (prior: Dirichlet)
 - Linear regression (prior: gaussian, doubly-exponential)
 -

Generalization via Risk Minimization

Risk Minimization

True Risk

Target performance measure

$$R(f) := \mathbb{E}(\ell(f(X), Y))$$

Classification

Probability of misclassification

$$P(f(X) \neq Y)$$

Empirical Risk

Performance on training data

$$\hat{R}_D(f) := \frac{1}{|D|} \sum_{i \in D} \ell(f(X_i), Y_i)$$

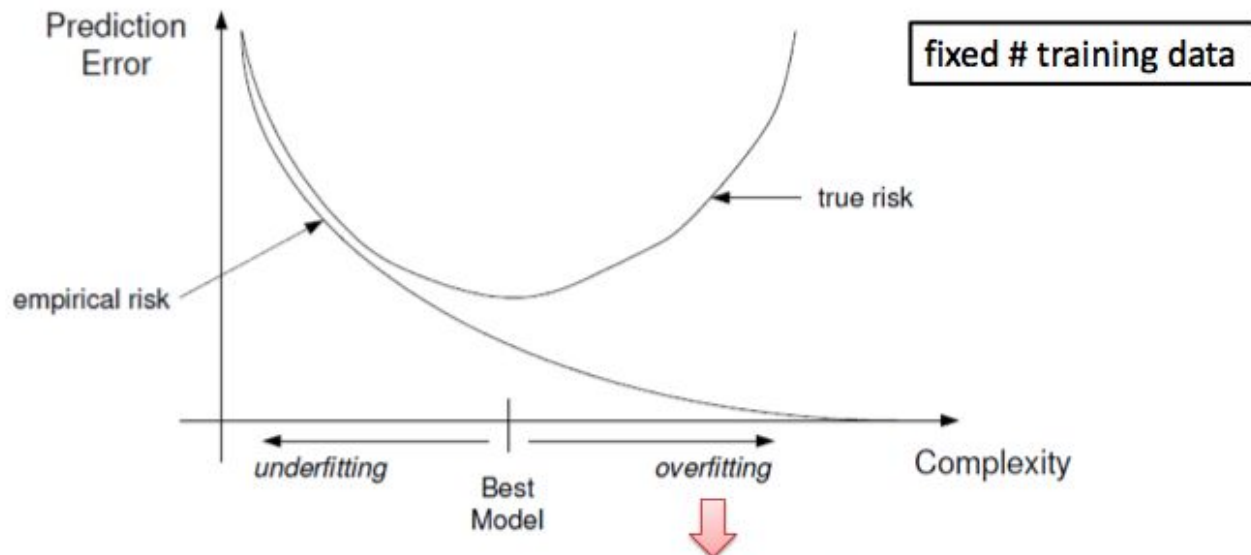
Classification

Proportion of misclassified examples

$$\frac{1}{n} \sum_{i=1}^n 1_{f(X_i) \neq Y_i}$$

Overfitting: Effect of discrepancy between empirical and true risks

If we allow very complicated predictors, we could overfit the training data.



Empirical risk is no longer a good indicator of true risk

Behavior of True Risk

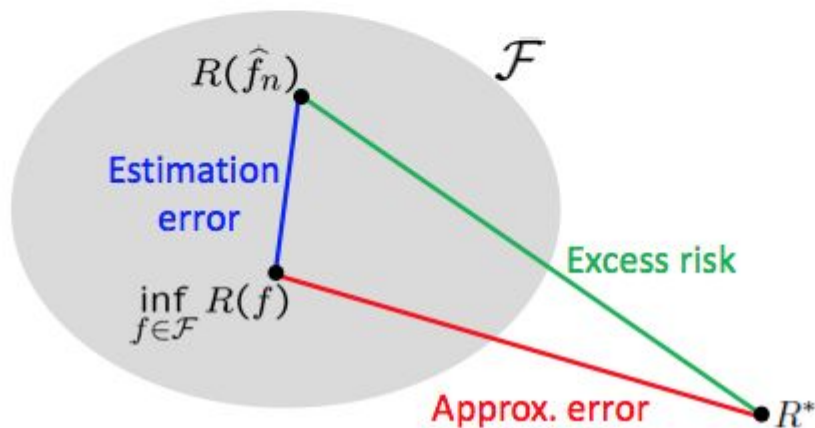
Want \hat{f}_n to be as good as optimal predictor f^*

Excess Risk $E[R(\hat{f}_n)] - R^* = \underbrace{\left(E[R(\hat{f}_n)] - \inf_{f \in \mathcal{F}} R(f)\right)}_{\text{estimation error}} + \underbrace{\left(\inf_{f \in \mathcal{F}} R(f) - R^*\right)}_{\text{approximation error}}$

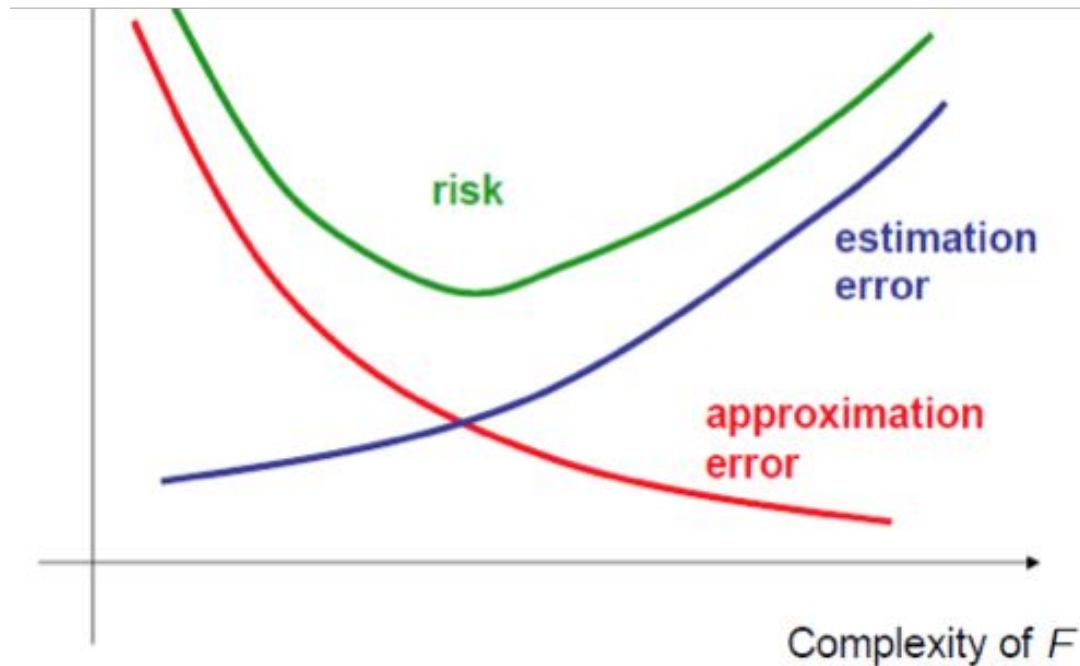
finite sample size
+ noise

Due to randomness
of training data

Due to restriction
of model class



Behavior of True Risk



fixed # training data

Linear Regression

- Assume $Y = f(X) + \epsilon$ where $\epsilon \sim \mathcal{N}(0, \sigma^2)$, and f does not need to be a linear function.
- Linear regression finds the best linear approximator:

$$\beta^* = \operatorname{argmin}_{\beta} \mathbb{E}[(Y - X\beta)^2] \quad \longleftarrow \text{Ideally}$$

$$\hat{\beta} = \operatorname{argmin}_{\beta} \frac{1}{n} \sum_{i=1}^n (Y_i - X_i\beta)^2 \quad \longleftarrow \text{Empirically}$$

where $X \in \mathbb{R}^{n \times p}$, $Y \in \mathbb{R}^n$, $\beta \in \mathbb{R}^p$.

- A polynomial regression “is” a linear regression. Instead, one would want to replace X with a polynomial basis (e.g. $\phi(X) = (X_1, X_2, X_1 X_2, X_1^2, \dots)$)

Linear Regression: Solving

- Fortunately, this is a convex optimization problem. Define

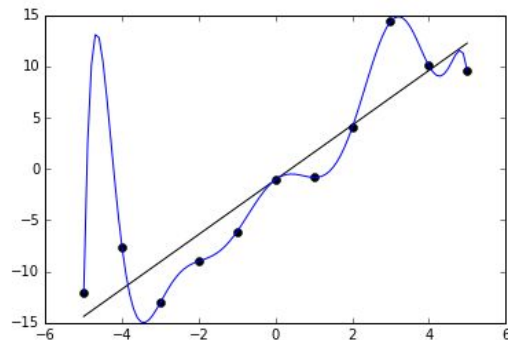
$$J(\beta) = (Y - X\beta)^T (Y - X\beta) = \sum_{i=1}^n (Y_i - X_i\beta)^2$$

Then the solution is given by $\nabla_{\beta} J(\beta) = \frac{\partial J}{\partial \beta} = 0$.

- Usually assume $X \in \mathbb{R}^{n \times p}$, $n > p$ (i.e. we have enough data). Why?
- $\hat{\beta} = (X^T X)^{-1} X^T Y$
- What if $X^T X$ is not invertible? We can constrain our answers by **regularization**.

Linear Regression: Overfitting & Regularization

- Let's say we are doing a polynomial regression



- One way to reduce overfitting is to regularize the model. In this case, in linear regression, we typically do this by adding **sparsity constraints** to β .
- How do we measure sparsity?

Linear Regression: Regularization

- A good measure to constrain how “large” β is is through L_p norm:

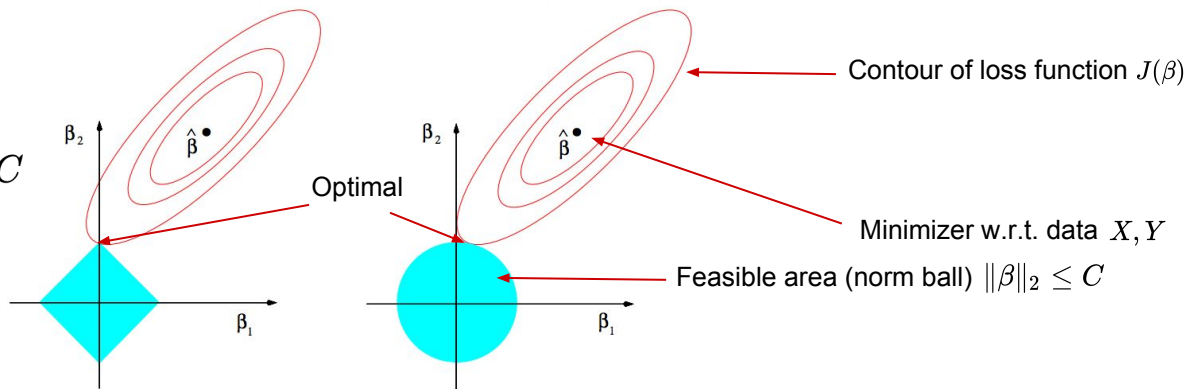
$$L_0 : \|\beta\|_0 = \sum_{i=1}^p \mathbf{1}_{\beta_i \neq 0}$$

$$L_1 : \|\beta\|_1 = \sum_{i=1}^p |\beta_i|$$

$$L_2 : \|\beta\|_2 = \sqrt{\sum_{i=1}^p \beta_i^2}$$

- Typically one would use L0, L1 (lasso), or L2 (ridge) regularizations.

$$\begin{aligned} & \min_{\beta} J(\beta) + \lambda \|\beta\|_p \\ \iff & \min_{\beta} J(\beta) \text{ s.t. } \|\beta\|_p \leq C \end{aligned}$$

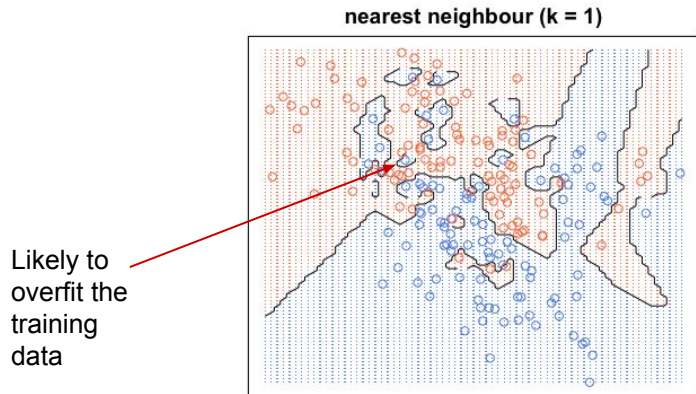


k Nearest Neighbors

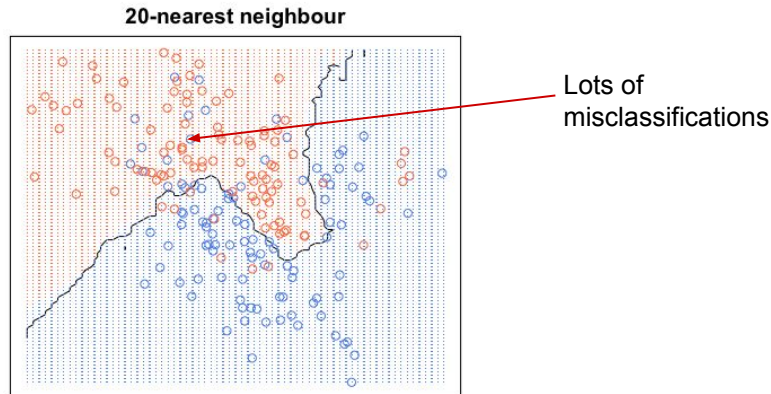
- High-level idea: predict at test time by **looking at the k nearest training data points**. More formally: $P(Y = c|X = x) = \frac{1}{k} \sum_{i: X_i \in \mathcal{N}_k(x)} 1_{Y_i=c}$.
- kNN **can** be used for regression as well: $\hat{f}(x) = \frac{1}{k} \sum_{i: X_i \in \mathcal{N}_k(x)} Y_i$.
- Efficient training: we only need to memorize all training data.
- Expensive evaluation: need to compute distance w.r.t. every training data point :-)
- **(T or F)** kNN algorithm works for **any** valid distance function/metric.
 - $d(x, y) \geq 0$, $d(x, y) = 0$ iff $x = y$, $d(x, y) = d(y, x)$, $d(x, z) \leq d(x, y) + d(y, z)$
Non-negativity Identity of indiscernibles Symmetry Triangular inequality

k Nearest Neighbors: Overfitting

- kNN can overfit. One way to correct this is to pick a better k using cross validation.



Low bias, high variance



High bias, low variance

k Nearest Neighbors: Properties & Summary

- **Non-parametric:** kNN does not assume the function form of the decision boundary. Helps us avoid the danger of mis-modeling the distribution.
- **Discriminative:** kNN models $P(Y|X)$, instead of $P(X, Y)$.
- **Instance-based:** kNN does not learn an explicit model. Instead, it classifies (or make prediction) based on training instances.
- **Curse of dimensionality:** kNN is bad at dealing with high dimensional data (where distance can be easily affected).

Naive Bayes

- Bayes Classifier with additional “naïve” assumption:

- Features are independent given class:

$$X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$$

$$\begin{aligned} P(X_1, X_2|Y) &= P(X_1|X_2, Y)P(X_2|Y) \\ &= P(X_1|Y)P(X_2|Y) \end{aligned}$$

- More generally:

$$P(X_1 \dots X_d|Y) = \prod_{i=1}^d P(X_i|Y)$$

$$X = \begin{bmatrix} X_1 \\ X_2 \\ \dots \\ X_d \end{bmatrix}$$

- If conditional independence assumption holds, NB is optimal classifier! But worse otherwise.

Naive Bayes

- Naive Bayes is a **generative** model with a prior on each class:

$$P(X, Y) = P(Y) \cdot \prod_{i=1}^n P(X_i | Y)$$

- At **training**, we need to find all parameters (see above). Typically found them via MLE estimate. For instance,

$$P_{MLE}(Y = c) = \frac{1}{n} \sum_{i=1}^n 1_{Y_i=c}$$

- At **evaluation**, we apply the Naive Bayes assumption on the argmax:

$$\hat{y} = \operatorname{argmax}_y P(Y = y | X) = \operatorname{argmax}_y P(X | Y = y) \cdot P(Y = y)$$

Logistic Regression

Assumes the following functional form for $P(Y|X)$:

$$P(Y = 0|X) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$\Rightarrow P(Y = 1|X) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$\Rightarrow \frac{P(Y = 1|X)}{P(Y = 0|X)} = \exp(w_0 + \sum_i w_i X_i) \begin{matrix} 1 \\ \geq \\ 0 \end{matrix} \geq 1$$

$$\Rightarrow w_0 + \sum_i w_i X_i \begin{matrix} 1 \\ \geq \\ 0 \end{matrix}$$

If >0 , then we
predict 1;
otherwise we
predict 0

Logistic Regression

Assumes the following functional form for $P(Y|X)$:

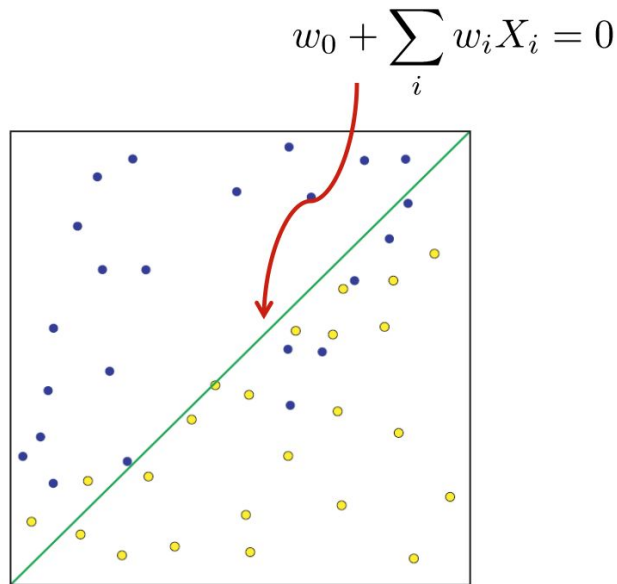
$$P(Y = 0|X) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

Decision boundary: Note - Labels are 0,1

$$P(Y = 0|X) \underset{1}{\overset{0}{\geq}} P(Y = 1|X)$$

$$w_0 + \sum_i w_i X_i \underset{0}{\overset{1}{\geq}} 0$$

(Linear Decision Boundary)



Logistic Regression

$$P(Y = 0|\mathbf{X}, \mathbf{w}) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$P(Y = 1|\mathbf{X}, \mathbf{w}) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$\begin{aligned} l(\mathbf{w}) &\equiv \ln \prod_j P(y^j | \mathbf{x}^j, \mathbf{w}) \\ &= \sum_j \left[y^j (w_0 + \sum_i^d w_i x_i^j) - \ln(1 + \exp(w_0 + \sum_i^d w_i x_i^j)) \right] \end{aligned}$$

Bad news: no closed-form solution to maximize $l(\mathbf{w})$

Good news: $l(\mathbf{w})$ is concave function of \mathbf{w}
concave functions easy to maximize

Logistic Regression

$$p(\mathbf{w} \mid Y, \mathbf{X}) \propto P(Y \mid \mathbf{X}, \mathbf{w})p(\mathbf{w})$$

- Define priors on \mathbf{w}
 - Common assumption: Normal distribution, zero mean, identity covariance
 - “Pushes” parameters towards zero

$$p(\mathbf{w}) = \prod_i \frac{1}{\kappa\sqrt{2\pi}} e^{\frac{-w_i^2}{2\kappa^2}}$$

Zero-mean Gaussian prior

- M(C)AP estimate $\mathbf{w}^* = \arg \max_{\mathbf{w}} \ln \left[p(\mathbf{w}) \prod_{j=1}^n P(y^j \mid \mathbf{x}^j, \mathbf{w}) \right]$

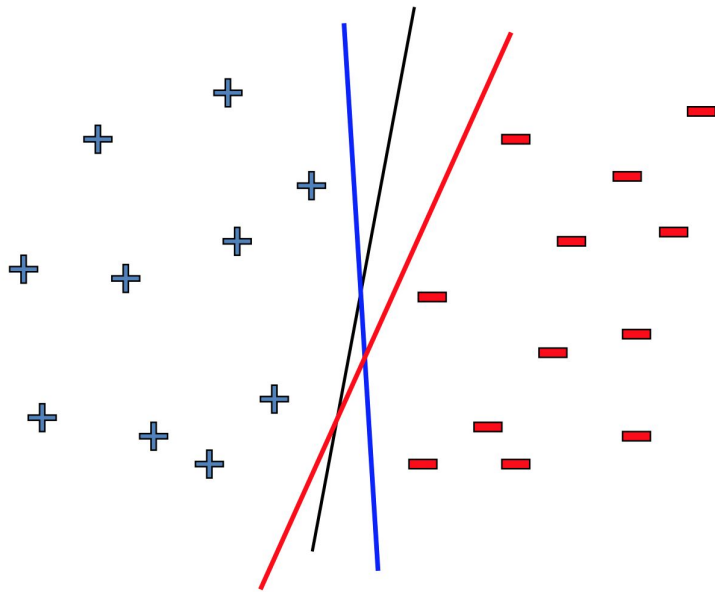
$$\mathbf{w}^* = \arg \max_{\mathbf{w}} \sum_{j=1}^n \ln P(y^j \mid \mathbf{x}^j, \mathbf{w}) - \underbrace{\sum_{i=1}^d \frac{w_i^2}{2\kappa^2}}$$

Still concave objective!

Penalizes large weights

Support Vector Machines (SVMs)

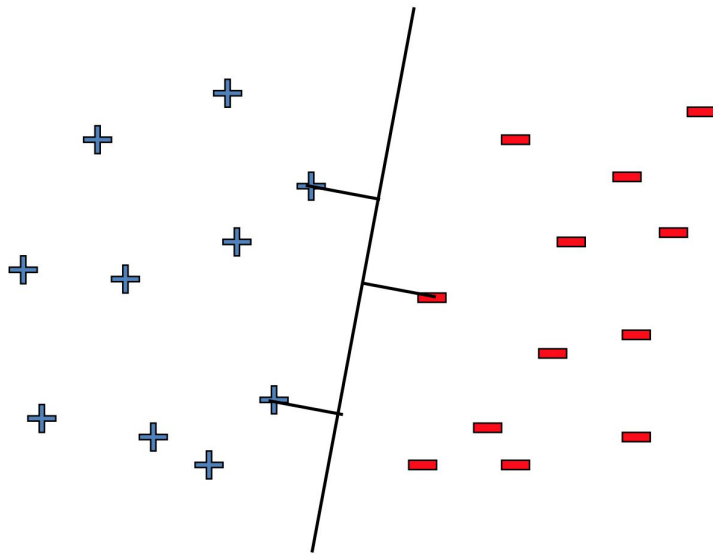
Consider a classification problem:



Often, the training data can be separated by multiple linear classifiers.
Which one to choose?

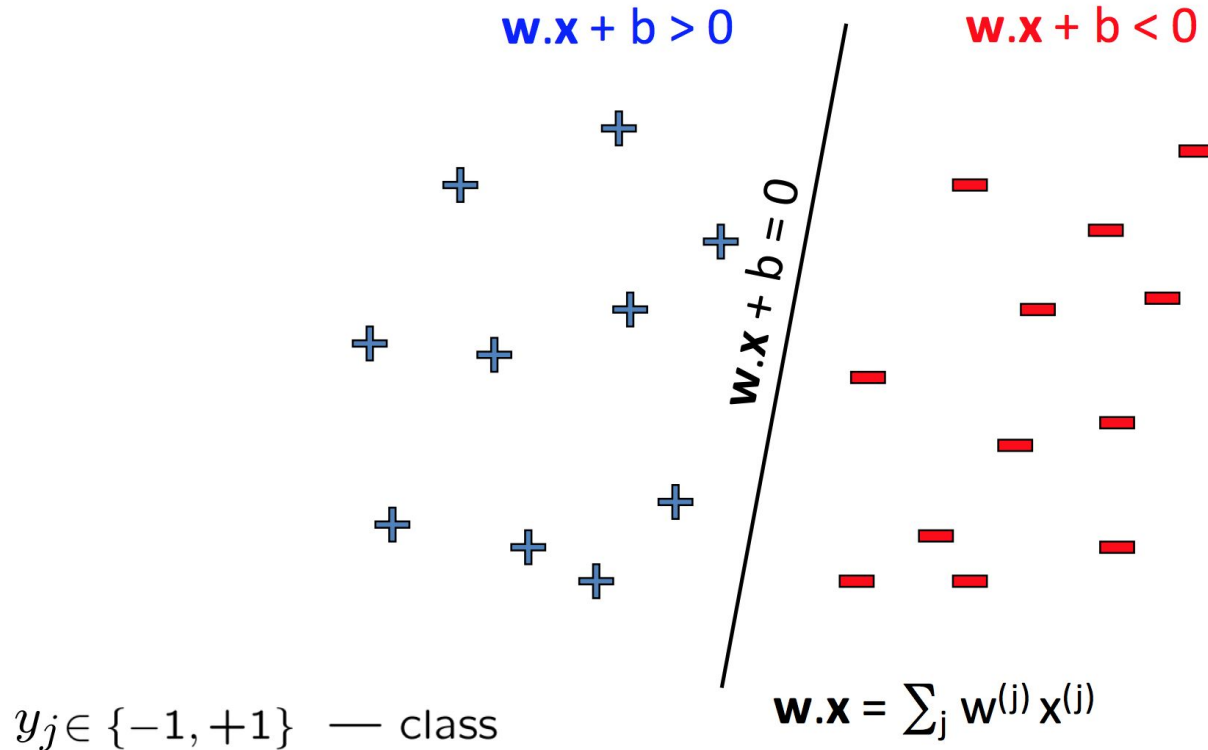
Support Vector Machines (SVMs)

Consider a classification problem:

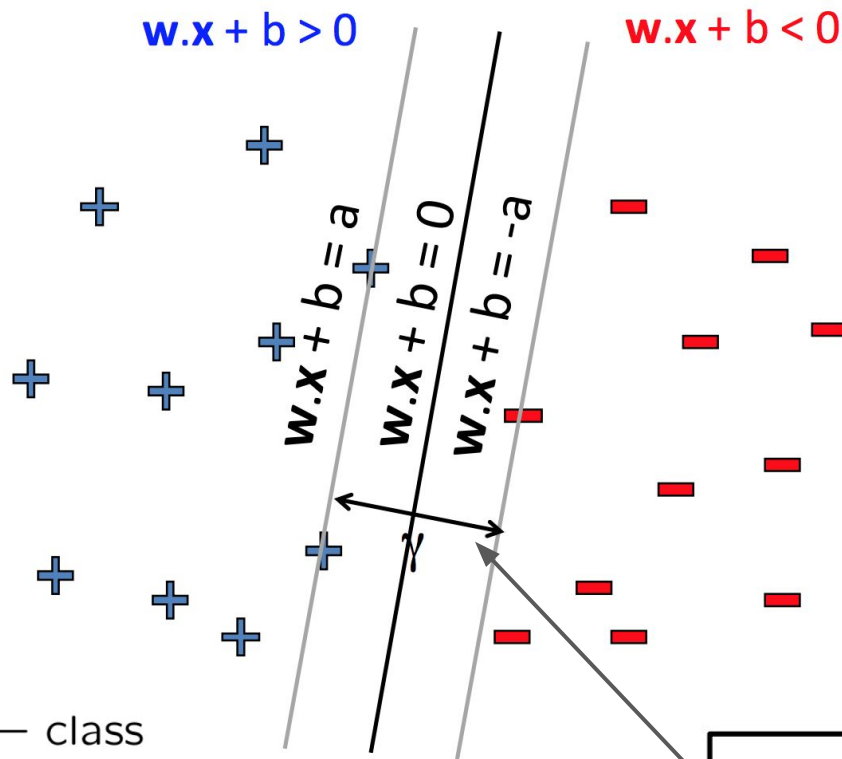


SVMs choose the one with the largest margin.

Support Vector Machines (SVMs)



Support Vector Machines (SVMs)



$y_j \in \{-1, +1\}$ — class

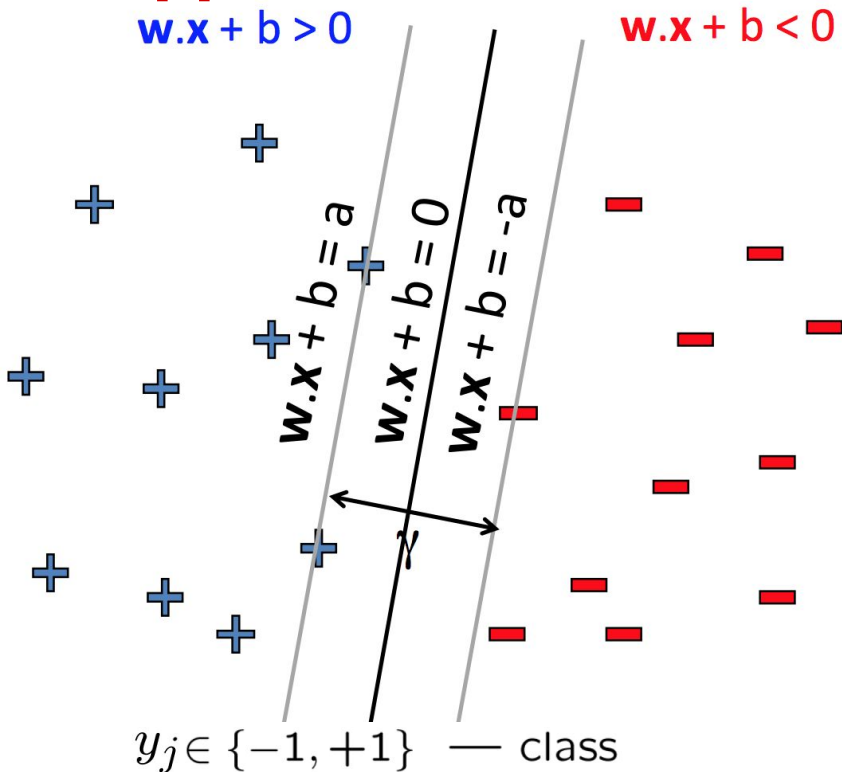
“confidence” $= (w \cdot x_j + b) y_j$

$$\text{margin} = \gamma = 2a/\|w\|$$

Support Vector Machines (SVMs)

$$w \cdot x + b > 0$$

$$w \cdot x + b < 0$$



Find w and b by solving:

$$\begin{aligned} \max_{w, b} \gamma &= \frac{2a}{\|w\|} \\ \text{s.t. } (w^T x_j + b)y_j &\geq a \quad \forall j \end{aligned}$$

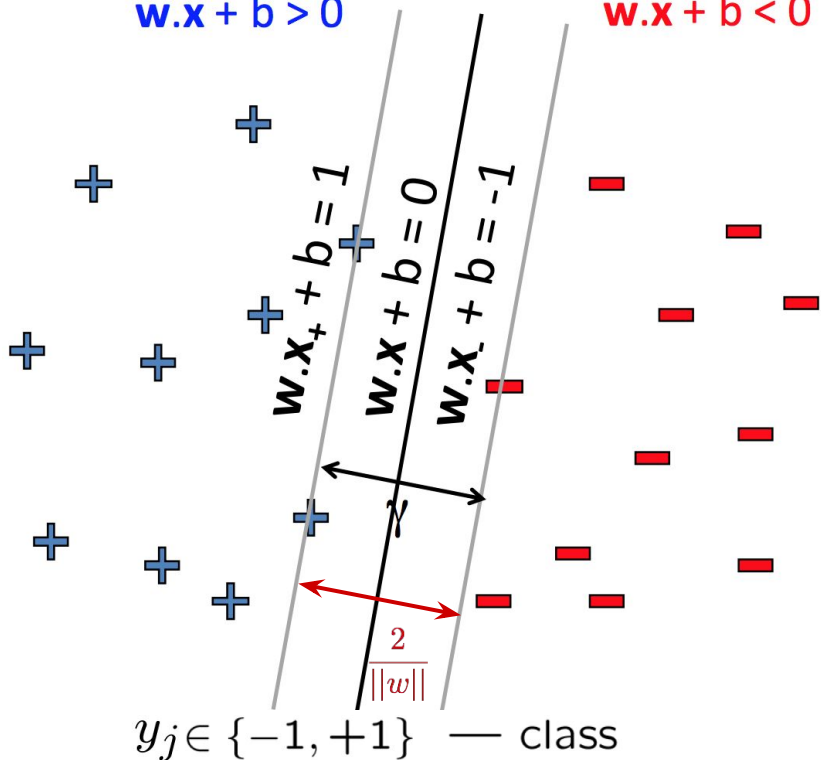
$$\text{margin} = \gamma = 2a/\|w\|$$

$$\text{"confidence"} = (w \cdot x_j + b) y_j$$

Support Vector Machines (SVMs)

$$w \cdot x + b > 0$$

$$w \cdot x + b < 0$$



Find w and b by solving:

$$\begin{aligned} \max_{w,b} \gamma &= \frac{2a}{\|w\|} \\ \text{s.t. } (w^T x_j + b)y_j &\geq a \quad \forall j \end{aligned}$$



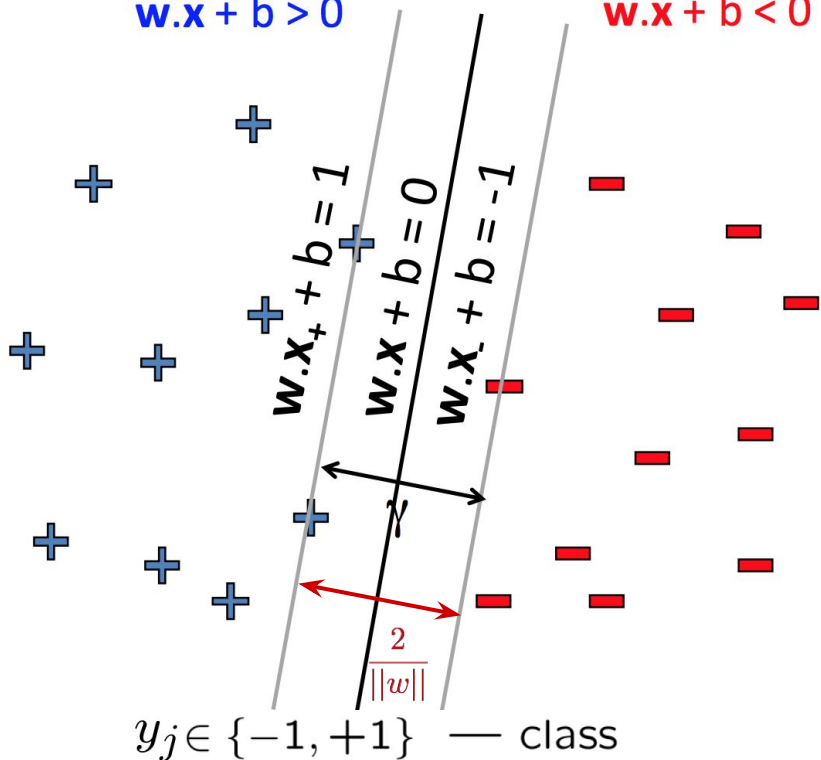
$$\begin{aligned} \max_{w,b} \frac{2}{\|w\|} \\ \text{s.t. } (w^T x_j + b)y_j &\geq 1 \quad \forall j \end{aligned}$$

"confidence" $= (w \cdot x_j + b) y_j$

Support Vector Machines (SVMs)

$$w \cdot x + b > 0$$

$$w \cdot x + b < 0$$



Find w and b by solving:

$$\begin{aligned} \max_{w,b} \gamma &= \frac{2a}{||w||} \\ \text{s.t. } (w^T x_j + b)y_j &\geq a \quad \forall j \end{aligned}$$



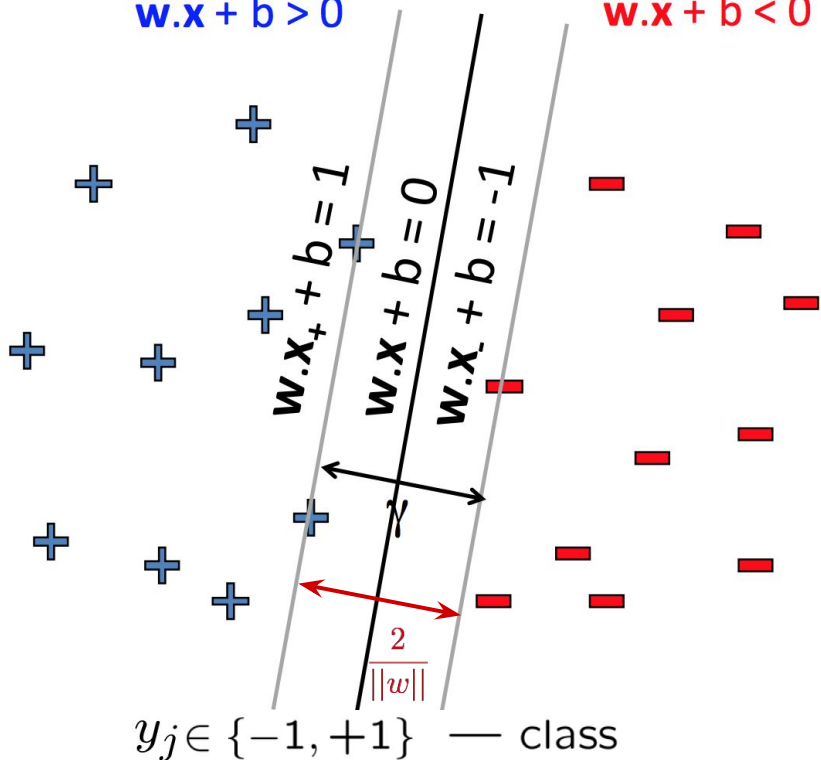
$$\begin{aligned} \min_{w,b} \frac{||w||}{2} \\ \text{s.t. } (w^T x_j + b)y_j &\geq 1 \quad \forall j \end{aligned}$$

“confidence” $= (w \cdot x_j + b) y_j$

Support Vector Machines (SVMs)

$$w \cdot x + b > 0$$

$$w \cdot x + b < 0$$



Find w and b by solving:

$$\begin{aligned} \max_{w,b} \gamma &= \frac{2a}{||w||} \\ \text{s.t. } (w^T x_j + b)y_j &\geq a \quad \forall j \end{aligned}$$



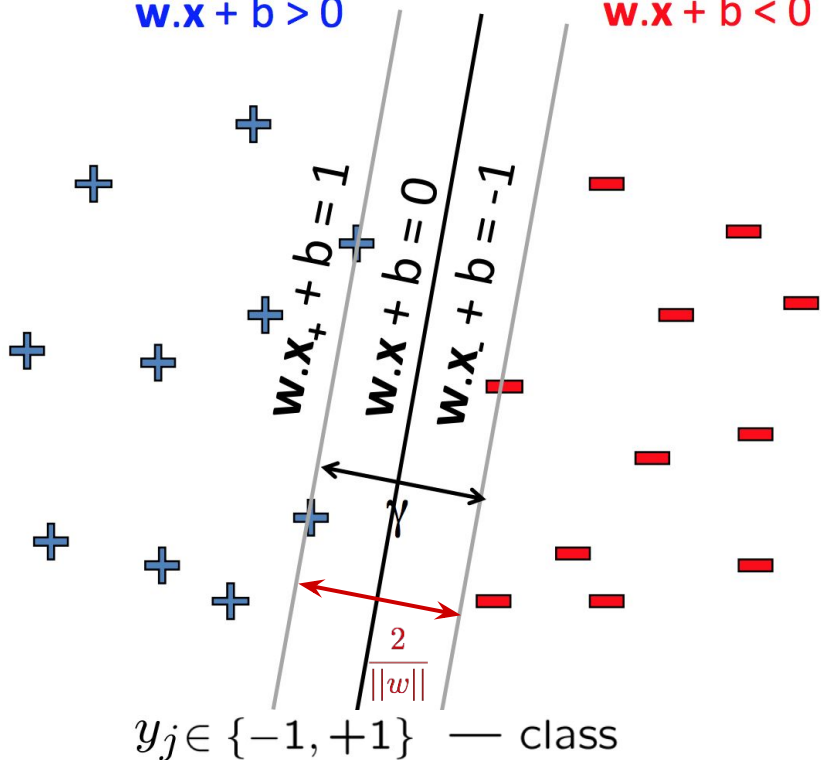
$$\begin{aligned} \min_{w,b} ||w|| \\ \text{s.t. } (w^T x_j + b)y_j &\geq 1 \quad \forall j \end{aligned}$$

"confidence" $= (w \cdot x_j + b) y_j$

Support Vector Machines (SVMs)

$$w \cdot x + b > 0$$

$$w \cdot x + b < 0$$



Find w and b by solving:

$$\begin{aligned} \max_{w,b} \gamma &= \frac{2a}{\|w\|} \\ \text{s.t. } (w^T x_j + b)y_j &\geq a \quad \forall j \end{aligned}$$



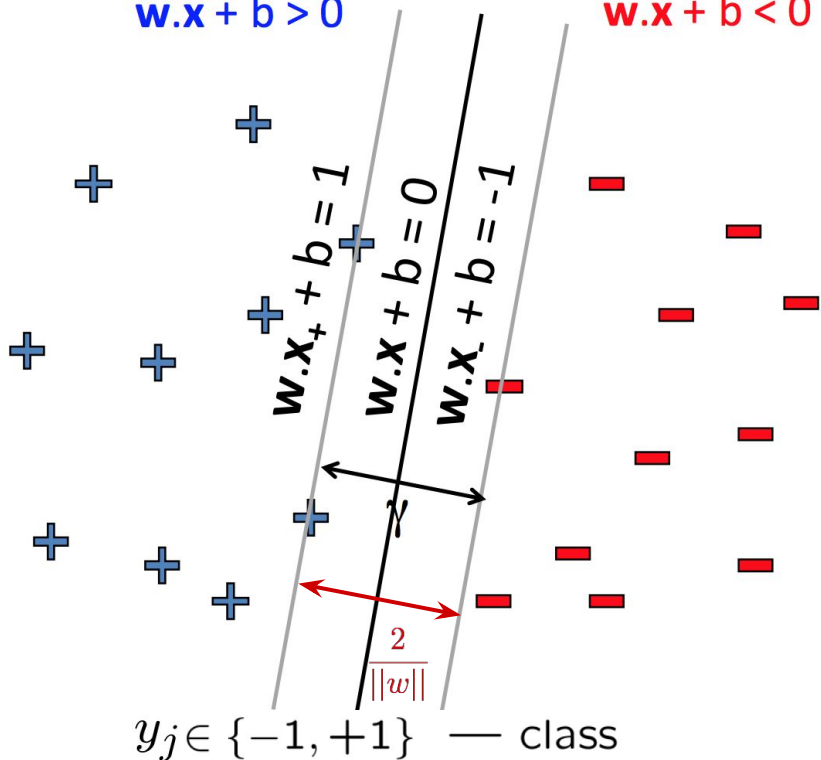
$$\begin{aligned} \min_{w,b} w^T w \\ \text{s.t. } (w^T x_j + b)y_j &\geq 1 \quad \forall j \end{aligned}$$

“confidence” $= (w \cdot x_j + b) y_j$

Support Vector Machines (SVMs)

$$w \cdot x + b > 0$$

$$w \cdot x + b < 0$$



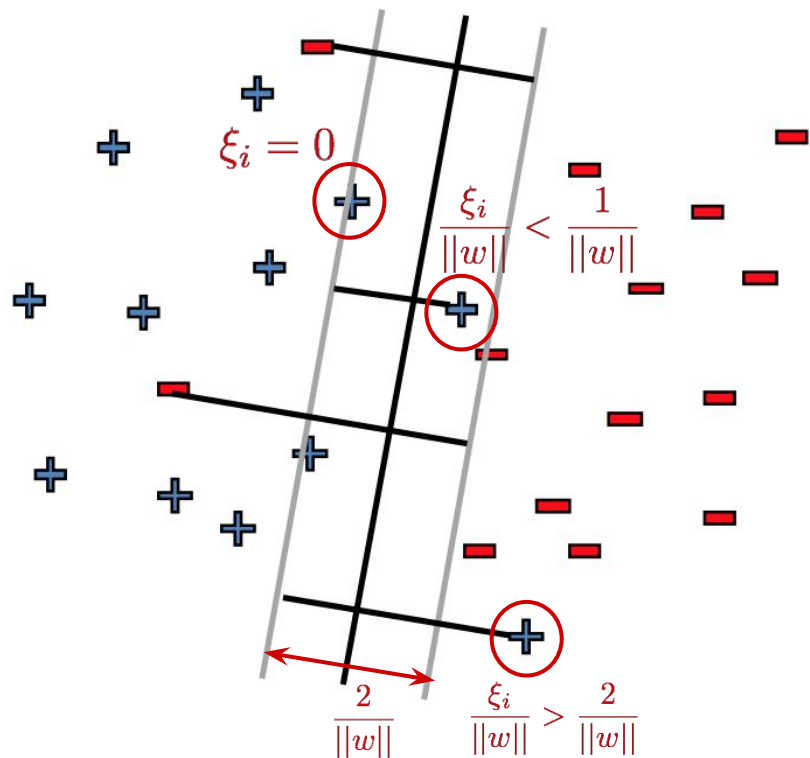
Find w and b by solving:

$$\begin{aligned} \min_{w, b} \quad & w^T w \\ \text{s.t.} \quad & (w^T x_j + b) y_j \geq 1 \quad \forall j \end{aligned}$$

We can solve this efficiently using quadratic programming (QP).

$$\text{"confidence"} = (w \cdot x_j + b) y_j$$

SVMs: What if data is not linearly separable?



Introduce “**slack**” variables!

$$\xi_i \geq 0$$

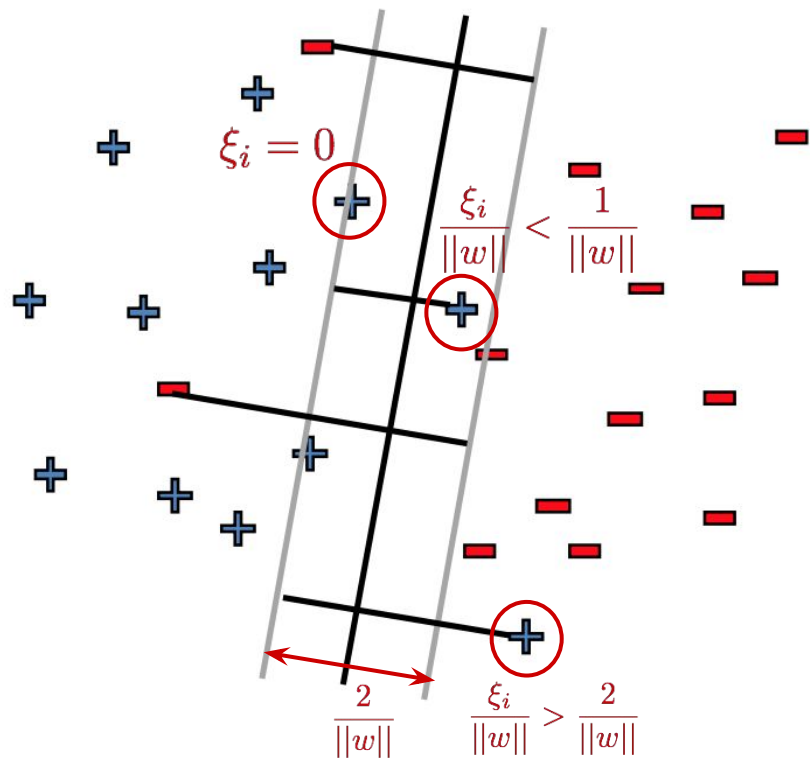
such that:

- $\xi_i = 0$ if the point is on the margin
- $0 < \xi_i \leq 1$ if point is between the margin and the correct side of the hyperplane
- $\xi_i > 1$ if the point is misclassified

The optimization problem becomes:

$$\begin{aligned} \min_{w, b, \xi_j} \quad & w^T w + C \sum_j \xi_j \\ \text{s.t.} \quad & (w^T x_j + b)y_j \geq 1 - \xi_j \quad \forall j \\ & \xi_j \geq 0 \quad \forall j \end{aligned}$$

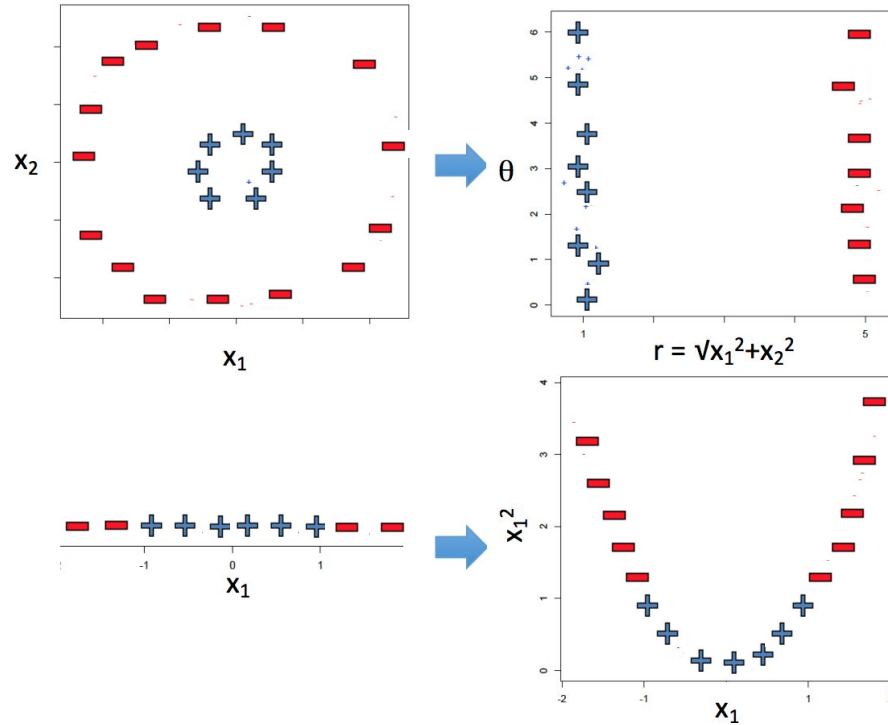
SVMs: What if data is not linearly separable?



$$\begin{aligned} \min_{w, b, \xi_j} \quad & w^T w + C \sum_j \xi_j \\ \text{s.t.} \quad & (w^T x_j + b)y_j \geq 1 - \xi_j \quad \forall j \\ & \xi_j \geq 0 \quad \forall j \end{aligned}$$

- Every constraint can be satisfied if ξ_i is sufficiently large.
- **C** is a **regularization** parameter:
 - C small allows constraints to be easily ignored \rightarrow large margin
 - C large makes constraints hard to ignore \rightarrow narrow margin
 - C = ∞ enforces all constraints \rightarrow hard margin
- This is still a **quadratic optimization problem** and there is a **unique minimum**.

SVMs: What if data is not linearly separable?



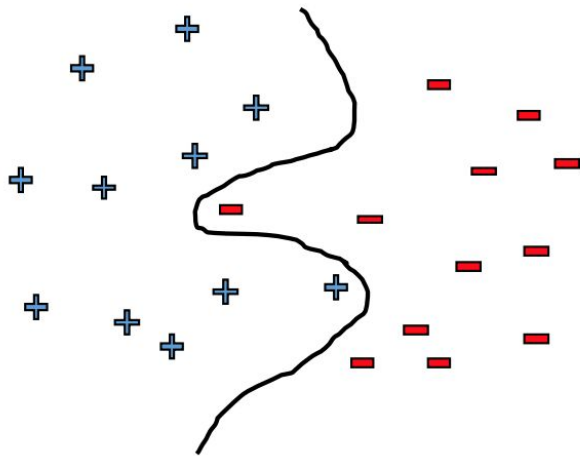
What if data is linearly separable using higher-order features ?

SVMs: What if data is not linearly separable?

Use a **feature mapping** $\Phi(x)$ to convert the data to a different space.

Example:

$$\Phi(\mathbf{x}) = (x_1^2, x_2^2, x_1x_2, \dots, \exp(x_1))$$



What if data is linearly separable using higher-order features ?

SVMs: What if data is not linearly separable?

In the original space:

$$\begin{aligned} & \min_{w,b} w^T w \\ \text{s.t. } & (w^T x_j + b)y_j \geq 1 \quad \forall j \end{aligned}$$

In feature space:

$$\begin{aligned} & \min_{w,b} w^T w \\ \text{s.t. } & (w^T \phi(x_j) + b)y_j \geq 1 \quad \forall j \end{aligned}$$

Use a **feature mapping** $\Phi(x)$ to convert the data to a different space.

Example:

$$\Phi(\mathbf{x}) = (x_1^2, x_2^2, x_1x_2, \dots, \exp(x_1))$$

Problem:

- $\Phi(x)$ can be very high dimensional (e.g. think of polynomial mapping).
- $\Phi(x)$ can be expensive to compute.

Solution:

- Maybe we don't need to compute $\Phi(x)$ explicitly!

SVMs: Deriving the dual problem

Primal problem:

$$\begin{aligned} \min_{w,b} \quad & w^T w \\ \text{s.t.} \quad & (w^T x_j + b)y_j \geq 1 \quad \forall j \end{aligned}$$

Introducing Lagrange multipliers:

$$L(w, b, \alpha) = \min_{w,b} \frac{1}{2} w^T w - \sum_j \alpha_j [(w^T x_j + b)y_j - 1]$$

Dual problem:

$$\begin{aligned} \max_{\alpha} \min_{w,b} \quad & L(w, b, \alpha) \\ \text{s.t.} \quad & \alpha_j \geq 0 \quad \forall j \end{aligned}$$

$$\text{s.t. } \alpha_j \geq 0 \quad \forall j$$

↑
weights on constraints
(one per training point)

SVMs: Deriving the dual problem

Primal problem:

$$\begin{aligned} \min_{w,b} \quad & w^T w \\ \text{s.t.} \quad & (w^T x_j + b)y_j \geq 1 \quad \forall j \end{aligned}$$

Introducing Lagrange multipliers:

$$L(w, b, \alpha) = \min_{w,b} \frac{1}{2} w^T w - \sum_j \alpha_j [(w^T x_j + b)y_j - 1]$$

Dual problem:

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i x_j \\ \text{s.t.} \quad & \sum_i \alpha_i y_i = 0 \quad \forall j \\ & \alpha_j \geq 0 \quad \forall j \end{aligned}$$

$$\text{s.t. } \alpha_j \geq 0 \quad \forall j$$

↑
weights on constraints
(one per training point)

SVMs: Deriving the dual problem

Primal problem:

$$\begin{aligned} \min_{w,b} \quad & w^T w \\ \text{s.t.} \quad & (w^T x_j + b)y_j \geq 1 \quad \forall j \end{aligned}$$



Primal in feature space:

$$\begin{aligned} \min_{w,b} \quad & w^T w \\ \text{s.t.} \quad & (w^T \phi(x_j) + b)y_j \geq 1 \quad \forall j \end{aligned}$$

Dual problem:

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i x_j \\ \text{s.t.} \quad & \sum_i \alpha_i y_i = 0 \quad \forall j \\ & \alpha_j \geq 0 \quad \forall j \end{aligned}$$




Dual in feature space:

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \underbrace{\phi(x_i) \phi(x_j)} \\ \text{s.t.} \quad & \sum_i \alpha_i y_i = 0 \quad \forall j \\ & \alpha_j \geq 0 \quad \forall j \end{aligned}$$



No need to compute $\Phi(x)$, only $\Phi(x)^T \Phi(x)$!

SVMs: Kernel trick

$$\begin{aligned} \max_{\alpha} \quad & \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \underbrace{\phi(x_i) \phi(x_j)} \\ \text{s.t.} \quad & \sum_i \alpha_i y_i = 0 \quad \forall j \\ & \alpha_j \geq 0 \quad \forall j \end{aligned}$$


No need to compute $\Phi(x)$, only $\Phi(x)^T \Phi(x)$!

Kernel trick: replace these dot products with an equivalent kernel function

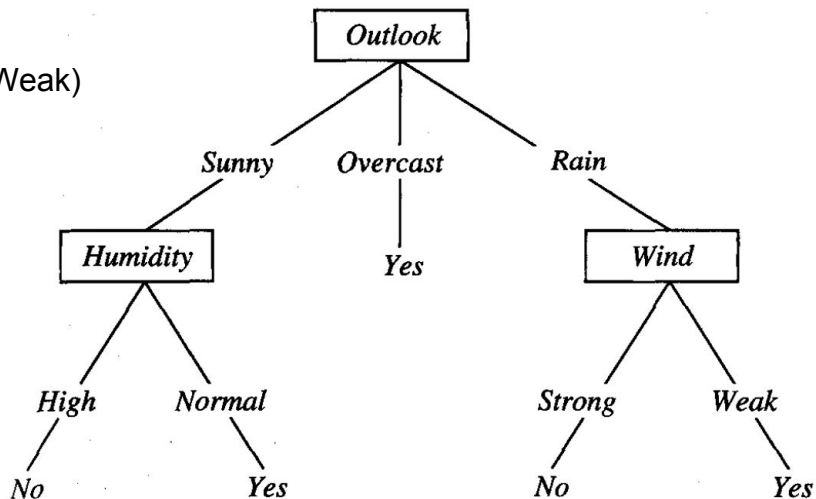
$$K(x_i, x_j) = \phi(x_i) \phi(x_j)$$

No need to know $\Phi(x)$, we can choose K directly from one of the commonly used kernel functions (e.g., polynomial, RBF, Gaussian, sigmoid).

Decision Tree Learning: Introduction

- Approximate discrete valued target functions where outcome is in form of decision tree
 - Sort observations by attributes from root to leaf

(Outlook = Sunny, Humidity = High, Wind = Weak)



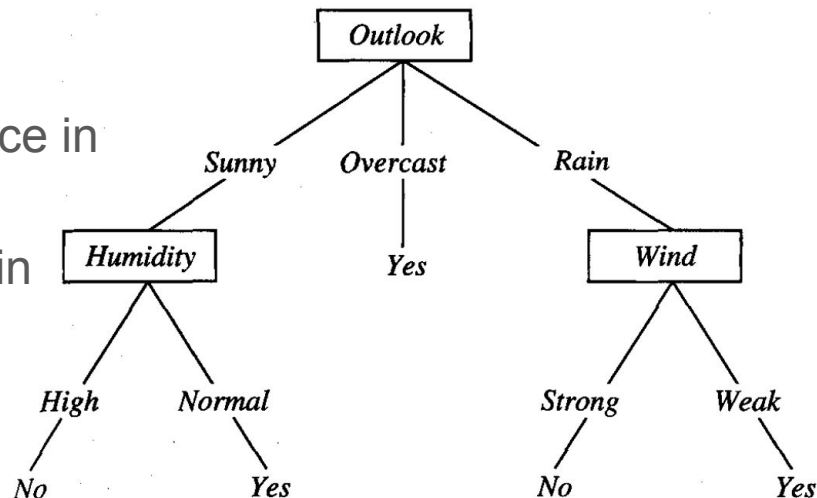
Decision Tree Learning: Usages

- Observations are represented by a fixed set of attribute pairs
- Target function has discrete action values
- Training Data may contain errors/is noisy/missing attribute values

Decision Tree Learning: Example Learner - ID3

- Learns decision trees in a top down manner
 - Sort instance attributes in descending order of importance from root to leaves.
 - Attribute order dependency contained within branch

- Each attribute can appear at most once in a branch
- Importance Measure: Information Gain



Decision Tree Learning: ID3 Importance Measure

Information Gain

- Advantage of attribute = decrease in uncertainty
 - Entropy of Y before split

$$H(Y) = - \sum_y P(Y = y) \log_2 P(Y = y)$$

- Entropy of Y after splitting based on X_i
 - Weight by probability of following each branch

$$\begin{aligned} H(Y | X_i) &= \sum_x P(X_i = x) H(Y | X_i = x) \\ &= - \sum_x P(X_i = x) \sum_y P(Y = y | X_i = x) \log_2 P(Y = y | X_i = x) \end{aligned}$$

- Information gain is difference

$$I(Y, X_i) = H(Y) - H(Y | X_i)$$

Max Information gain = min conditional entropy

Decision Tree Learning: ID3 Importance Measure

Pick the attribute/feature which yields maximum information gain:

$$\begin{aligned}\arg \max_i I(Y, X_i) &= \arg \max_i [H(Y) - H(Y|X_i)] \\ &= \arg \min_i H(Y|X_i)\end{aligned}$$

Entropy of Y

$$H(Y) = - \sum_y P(Y = y) \log_2 P(Y = y)$$

Conditional entropy of Y

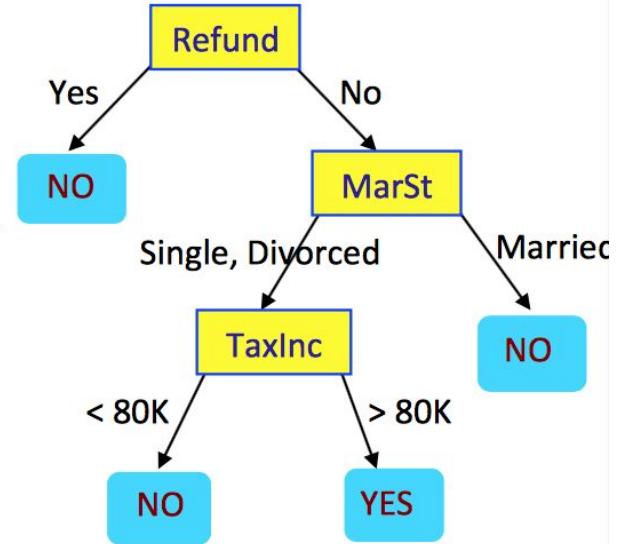
$$H(Y | X_i) = \sum_x P(X_i = x) H(Y | X_i = x)$$

Feature which yields maximum reduction in entropy (uncertainty)
provides maximum information about Y

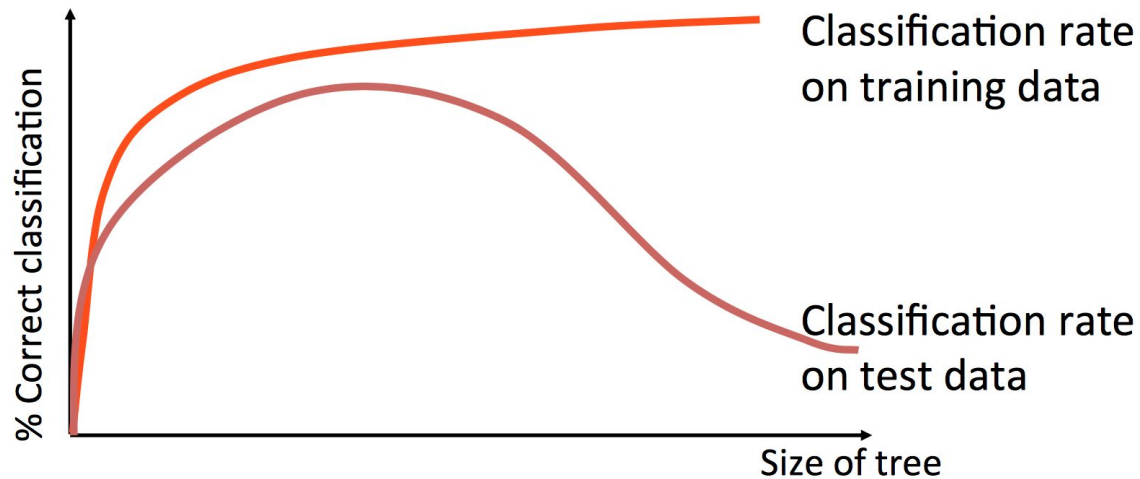
Decision Tree Learning: ID3 Build-Up

Main loop: C4.5

1. $X \leftarrow$ the “best” decision feature for next *node*
2. Assign X as decision feature for *node*
3. For “best” split of X , create new descendants of *node*
4. Sort training examples to leaf nodes
5. If training examples perfectly classified, Then STOP, Else iterate over new leaf nodes



Decision Tree Learning: Overfitting



- Can be reduced by pruning

Decision Tree Learning: General Takeaways

- In general a decision tree can represent any discrete valued function
- Decision Trees perform a complete search through the Hypothesis space (consisting of decision trees)
 - The target optimal decision tree is guaranteed to be in this space
- ID3 Algorithm creates trees by placing attributes with highest IG closest to root

Neural Networks

Consider a prediction problem in which we want to predict some **outputs** **y** from some **inputs** **x**.

The **goal** is to learn a **function** **f** such that

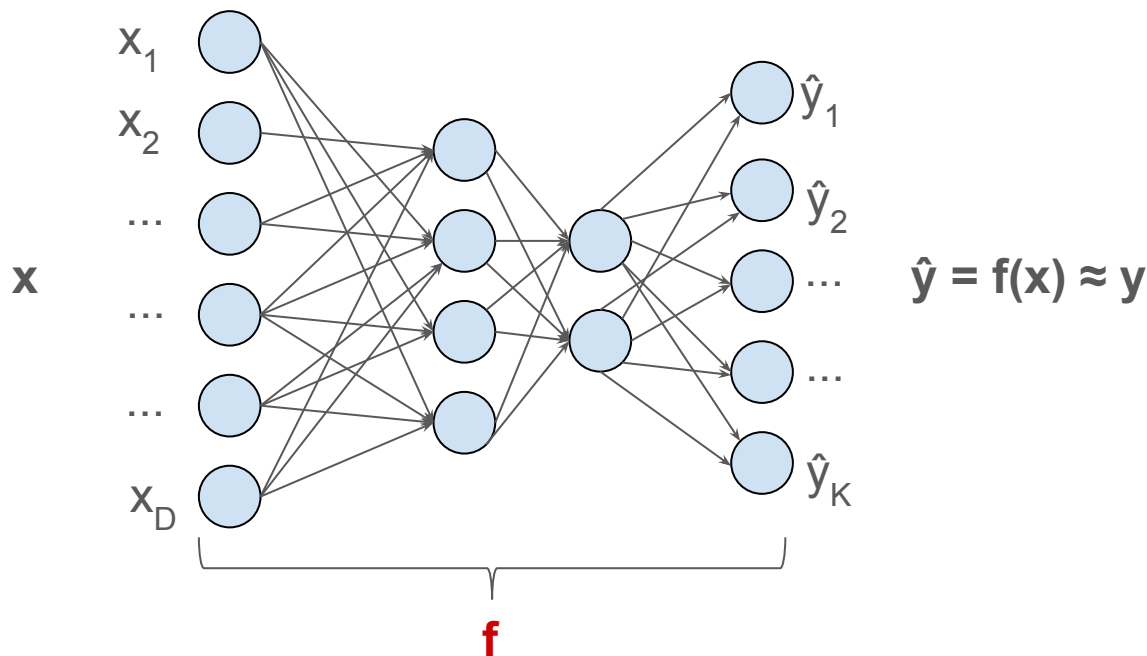
$$\mathbf{f}(\mathbf{x}) = \mathbf{y}$$

Examples of such functions learnt in class:

- Linear regression: $f(x) = w^T x + b$
- Logistic regression: $f(x) = P(y|x) = \frac{\exp(w^T x + b)}{1 + \exp(w^T x + b)}$

Neural Networks

A neural network is also such a function **f**, whose parameters are the weights of the network. For example:



Neural Networks: Artificial neuron

Neuron pre-activation:

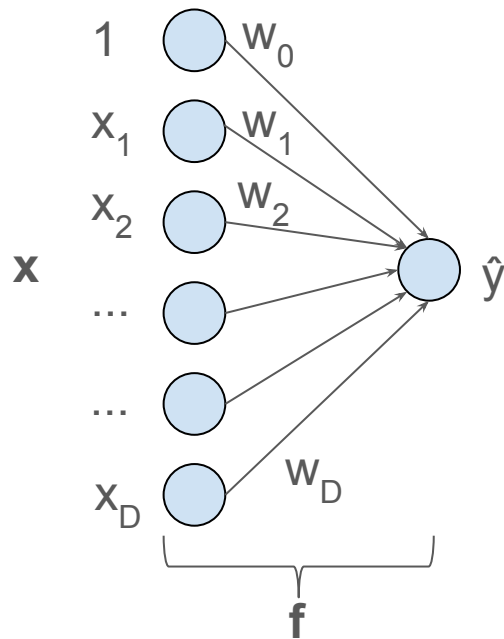
$$g(x) = w_0 + \sum_i w_i x_i = w_0 + w^T x$$

↑ ↑
bias weights

Neuron post-activation:

$$\hat{y} = a(g(x)) = a(w_0 + w^T x)$$

↑
activation
function



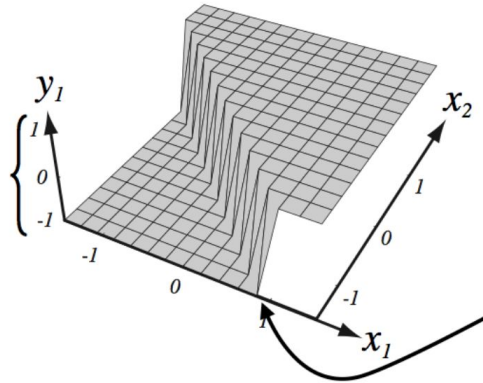
Neural Networks: Artificial neuron

Neuron post-activation:

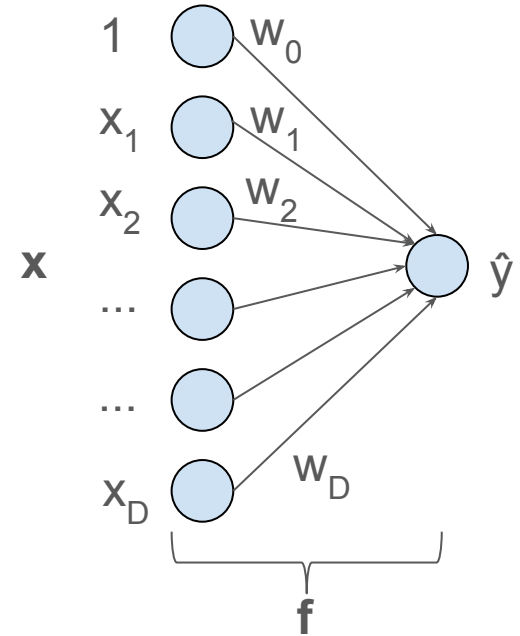
$$\hat{y} = a(g(x)) = a(w_0 + w^T x)$$

activation
function

The range of y is
determined by the
function $a(\cdot)$

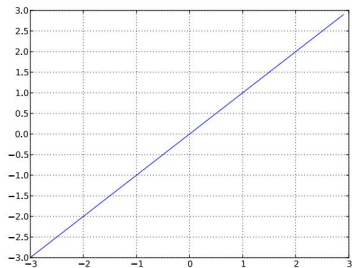


The bias changes
only the position

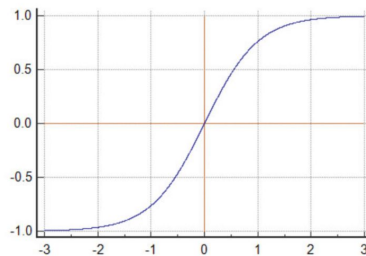


Neural Networks: Activation functions

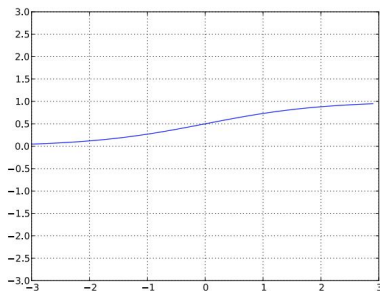
- **Linear:** $a(x) = c \cdot x + b$



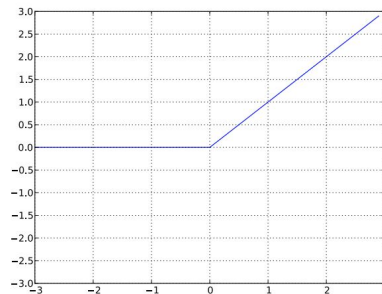
- **Tanh:** $a(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$



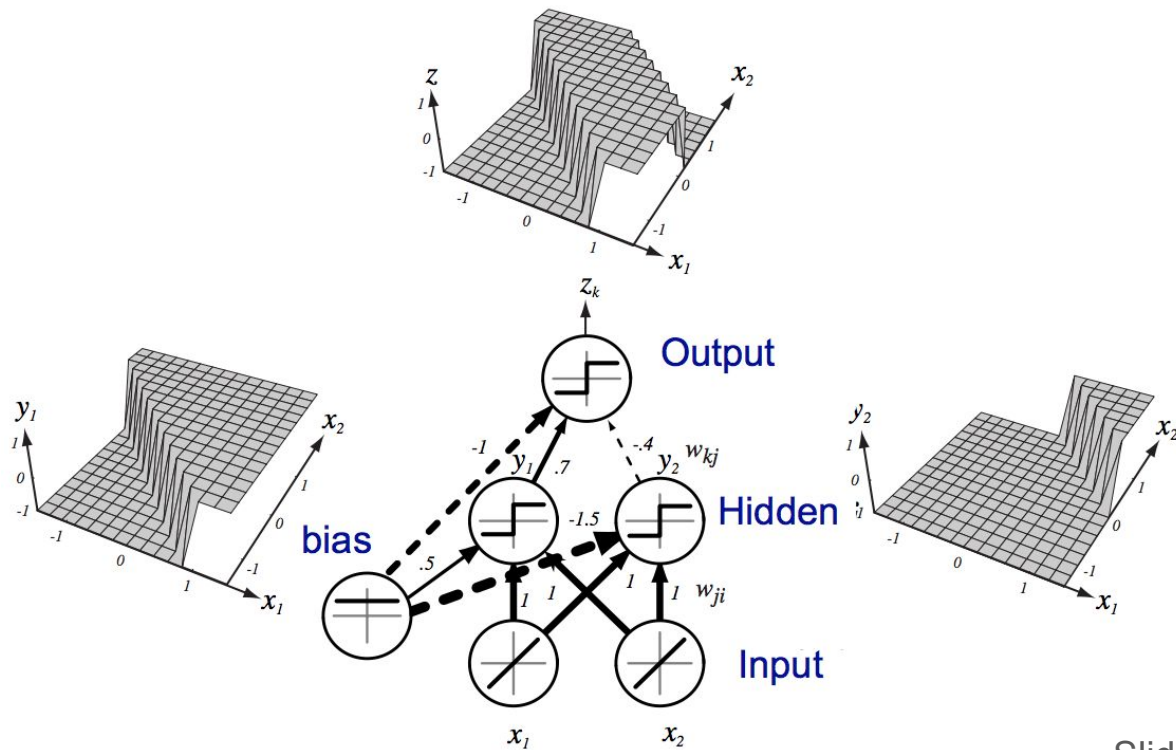
- **Sigmoid:** $a(x) = \frac{1}{1 + \exp(-x)}$



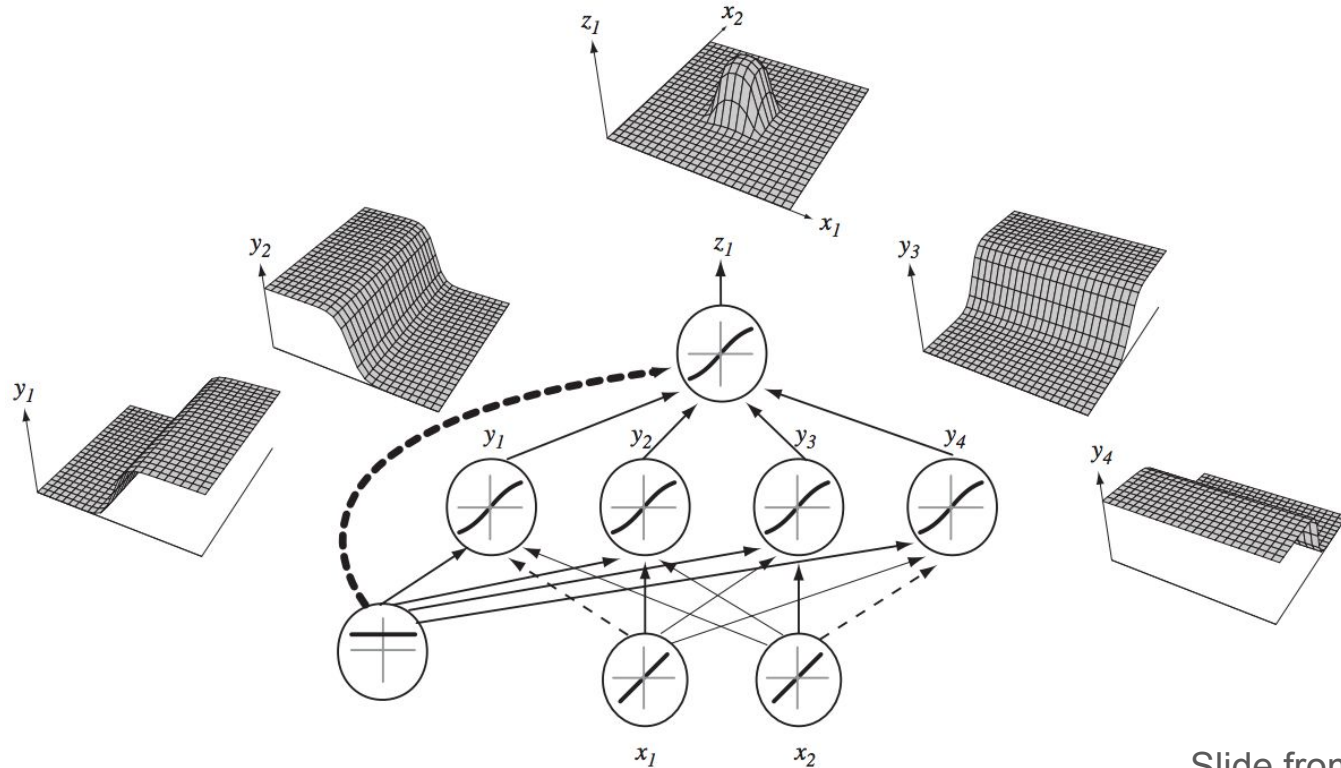
- **ReLU:** $a(x) = \max(0, x)$



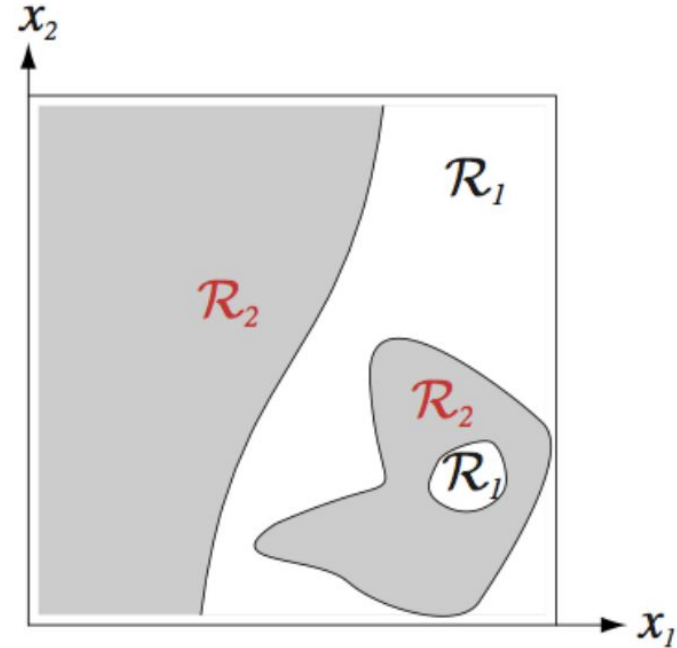
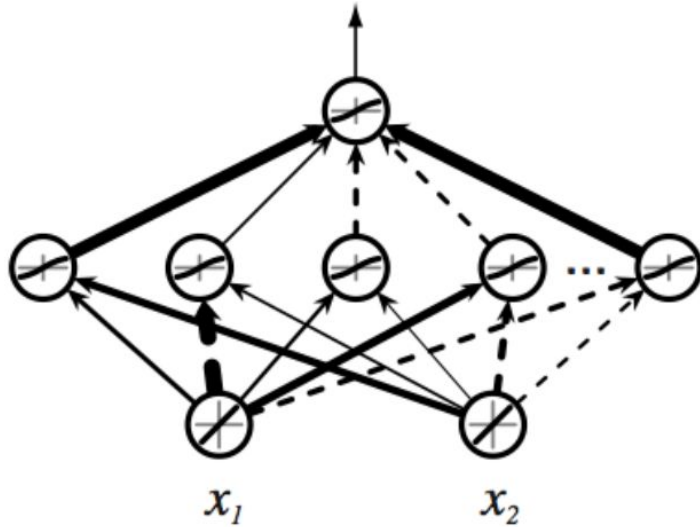
Neural Networks: Multiple layers



Neural Networks: Multiple layers



Neural Networks: Multiple layers



Neural Networks: Universal approximation

Universal Approximation Theorem (Hornik, 1991):

A single hidden layer neural network with a linear output unit can approximate any continuous function arbitrarily well, given enough hidden units.

Neural Networks: Universal approximation

Universal Approximation Theorem (Hornik, 1991):

A single hidden layer neural network with a linear output unit can approximate any continuous function arbitrarily well, given enough hidden units.

This does not mean that there is learning algorithm that can find the right parameter values!

Neural Networks: Training

- Learning is posed as an **optimization** problem.
- The goal is to minimize some **loss** function:

$$\min_{\theta=\{W,b\}} \underbrace{\ell(\hat{y}_{\text{pred}}, y)}_{=f(x)}$$

- One approach we learnt is to compute the **gradients** of the loss wrt. model parameters, and use them to update the parameters accordingly.
- **Gradient descent**: compute the gradient given N samples and update the parameters:

$$\theta^+ \leftarrow \theta - \alpha \nabla_{\theta} \sum_{i=1}^N \ell(\hat{y}(x^{(i)}), y^{(i)})$$

Neural Networks: Training

- How do we compute the gradients?

Backpropagation!

Useful resource:

<http://colah.github.io/posts/2015-08-Backprop/>

- Neural network training involves two phases:

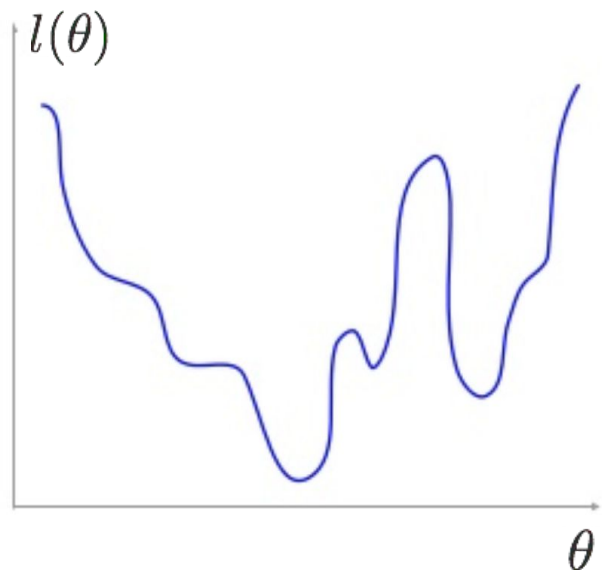
1. **Forward pass:**

Compute the function value $f(x)$ given some input x , and our current estimate of the network parameters.

2. **Backward pass (backpropagation):**

Compute all gradients $\nabla_{\theta} \ell$ with just one pass and update the parameters.

Neural Networks: Why we can't always find the optimal parameters



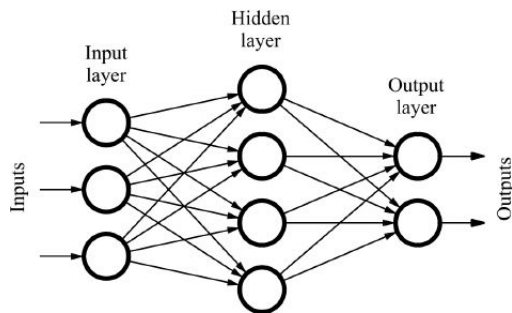
The loss function is usually not convex.

→ gradient descent will converge to a **local** optimum!

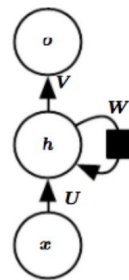
But that's ok!

Neural Networks: Examples

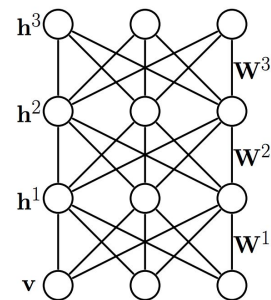
Fully-connected



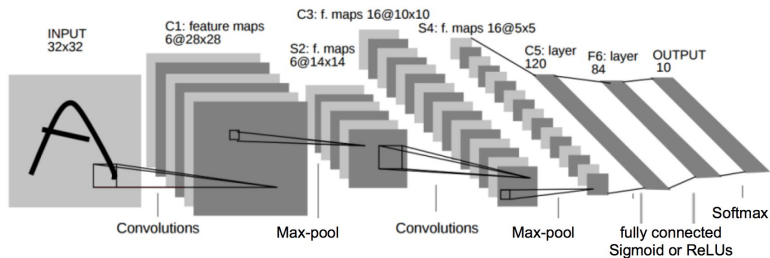
Recurrent



Deep Boltzmann Machine



Convolutional



Acknowledgements

- Neural network slides inspired by Russ Salakhutdinov's [slides](#)