

Machine Learning Overview

Yale

CPSC/AMTH/CBB 663

CPSC 452





THIS IS YOUR MACHINE LEARNING SYSTEM?

| YUP! YOU POUR THE DATA INTO THIS BIG
PILE OF LINEAR ALGEBRA, THEN COLLECT
THE ANSWERS ON THE OTHER SIDE.

| WHAT IF THE ANSWERS ARE WRONG?

| JUST STIR THE PILE UNTIL
THEY START LOOKING RIGHT.





Announcements

- Tomorrow at 3pm: Kincaid tutorial on Getting Set Up
- Sunday 9-11am: Creating a neural network in pytorch
 - See blog: <https://towardsdatascience.com/tabnet-deep-neural-network-for-structured-tabular-data-39eb4b27a9e4>
- Office hours begin next week!
- Homework 1 will be assigned TONIGHT, due Feb 15
- If you have questions to ask me, please ask via discord
 - → JOIN DISCORD!!
- Class has a partition on the HPC Server Farnam (see canvas announcement)



what does it mean for a machine or a program learn?

- ① minimize error on unseen data
- ② tune params to a model to better fit the data
- ③ some set of params are updated to improve results of a target function
- ④ decode the features of a dataset
- ⑤ find latent pattern in data without explicit instruction
- ⑥ minimize cost function doesn't mean you are learning
- ⑦ teach itself to perform better
- ⑧ find pattern or make prediction for data

• Machine learning algorithm = algorithm that can learn from data
 ≠ meta-learning

• What do we mean by “learn from data”?

• “A computer program is said to learn from : canonical def ML

- experience E
- with respect to some class of tasks T
- performance measure P

• IF its performance at tasks in T , as measured by P , improves with experience E .” – Mitchell (1997)

• What are E , T , and P ?

- It depends on the problem!

some abstract problems are hard to find P

data: text, image, numerical, tabular

task: human can do

P : same data + different $P \rightarrow$ different solutions

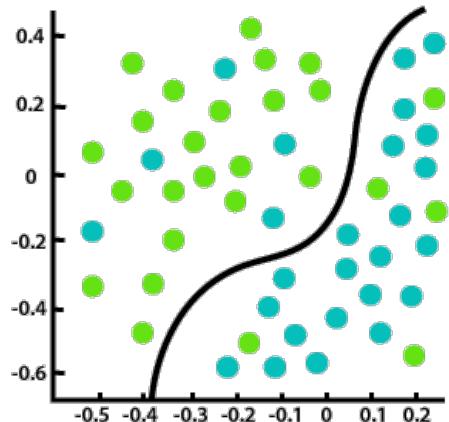


- Machine learning helps us to solve problems that are too difficult for fixed programs
- Learning is the process of attaining the ability to perform the task
- Tasks are typically described in terms of how to process an ***example***
 - An example is a collection of measured ***features***
 - Often represented as a vector $x \in \mathbb{R}^d$ of features
 - E.g. an example includes a single image with the features comprised of the pixels in the image

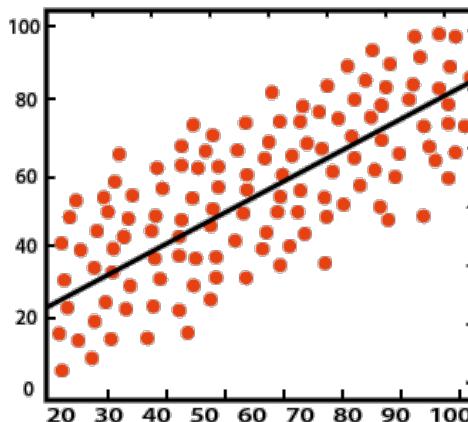
often in high dim



- Classification *draw a boundary in a feature space*
 - Output: a function $f: \mathbb{R}^d \rightarrow \{1, \dots, k\}$
- Regression
 - Output: a function $f: \mathbb{R}^d \rightarrow \mathbb{R}$



Classification



Regression

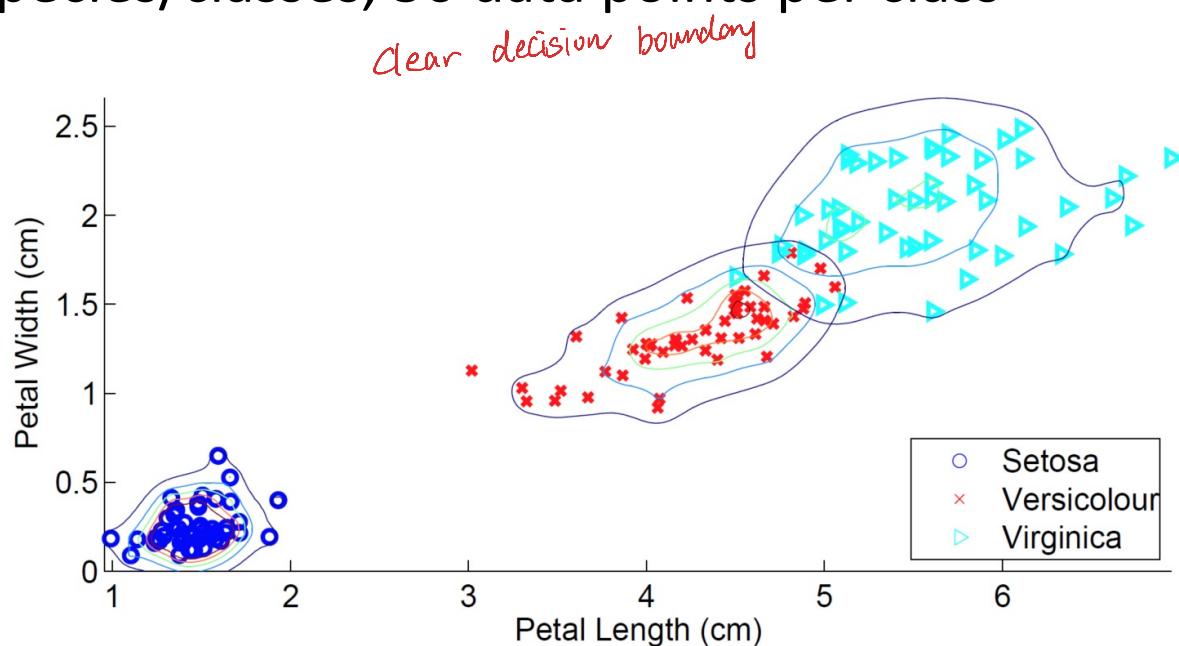
Linear regression example



- Need a quantitative measure to evaluate a machine learning algorithm
- Provides a measure to guide the selection of an algorithm
- Examples:
 - **Accuracy** (or equivalently, the ***error rate***) in a classification problem
 - Prediction error
- Performance is typically ***estimated***/evaluated based on a ***test set*** of data (different from training data)
- Choosing a performance measure isn't always straightforward



- Most algorithms **experience** an entire **dataset** consisting of many examples or ***data points***
- Example: Iris data set (Fisher, 1936)
 - Four features: sepal length, sepal width, petal length, petal width
 - 3 species/classes, 50 data points per class





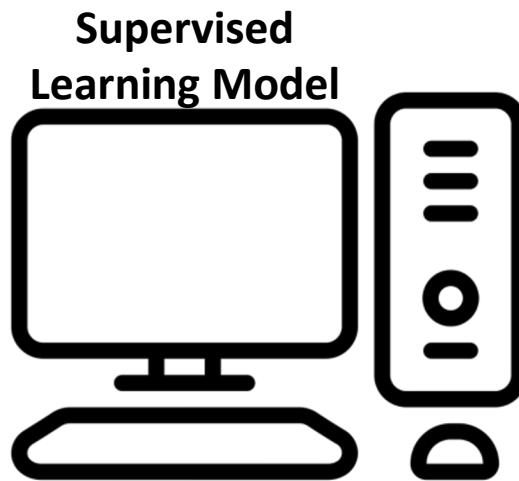
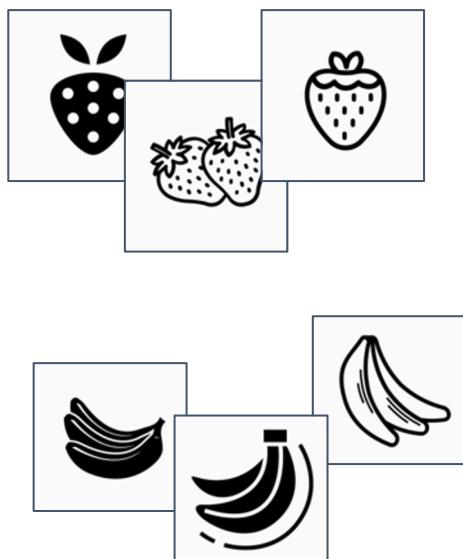
- Goal is to build a program that learns
 - from experience E
 - to do well at task T
 - as measured by performance measure P
- Correctly predicting the right classes
 - How would this be a performance measure?
 - Is it differentiable?

*we can't use accuracy because it jumps
often cross-entropy loss*

Supervised Learning

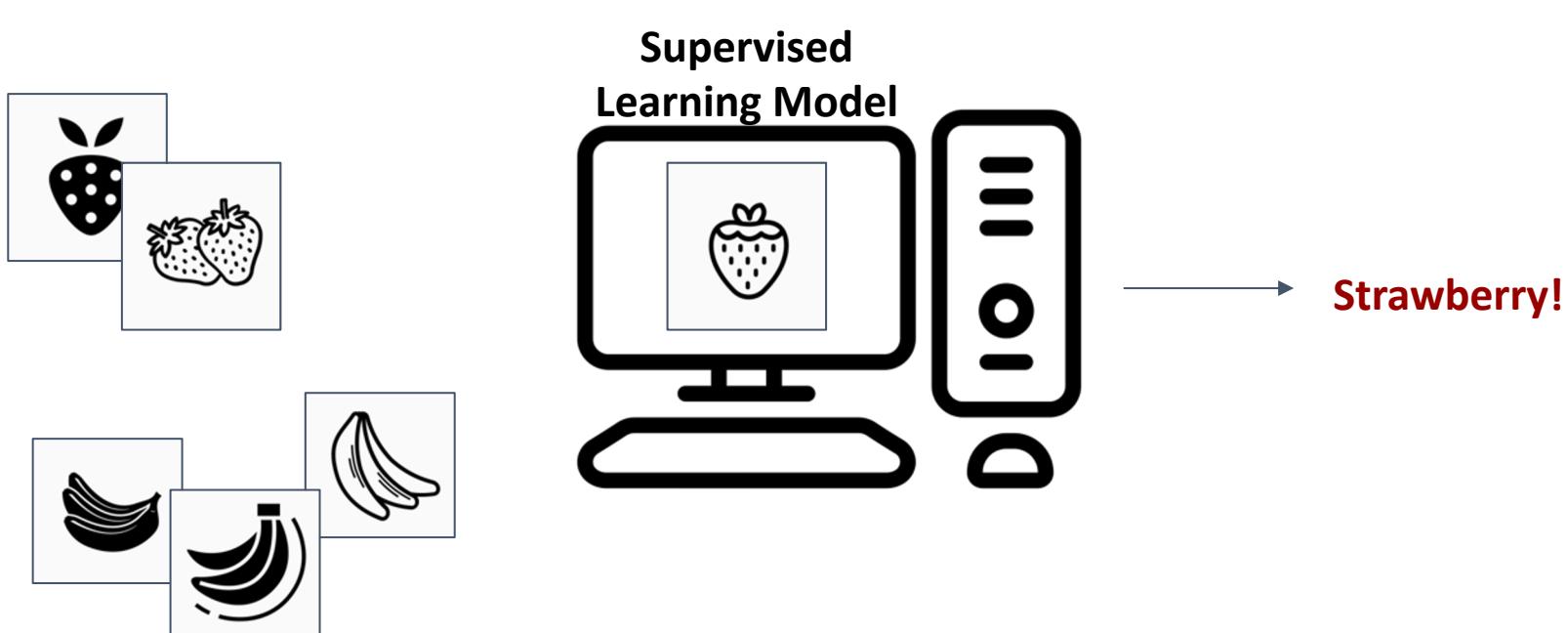
Supervised learning

- Have a bunch of labelled data,
want to label new data



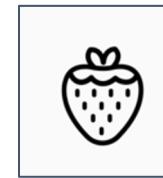
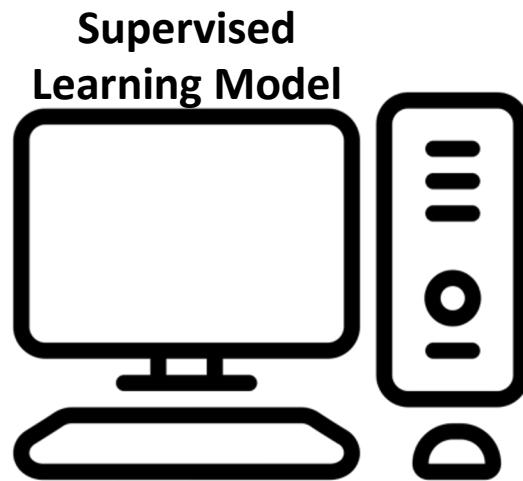
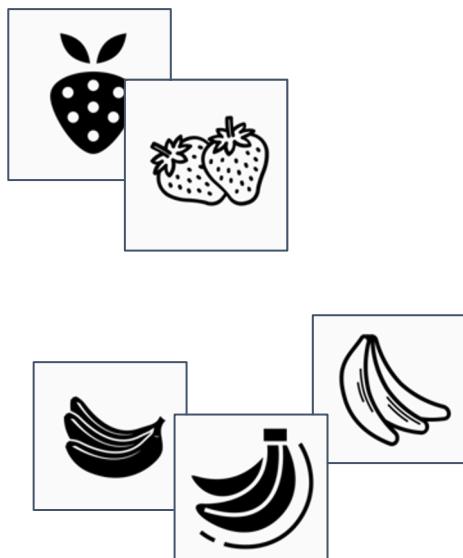
Supervised learning

- Have a bunch of labelled data,
want to label new data



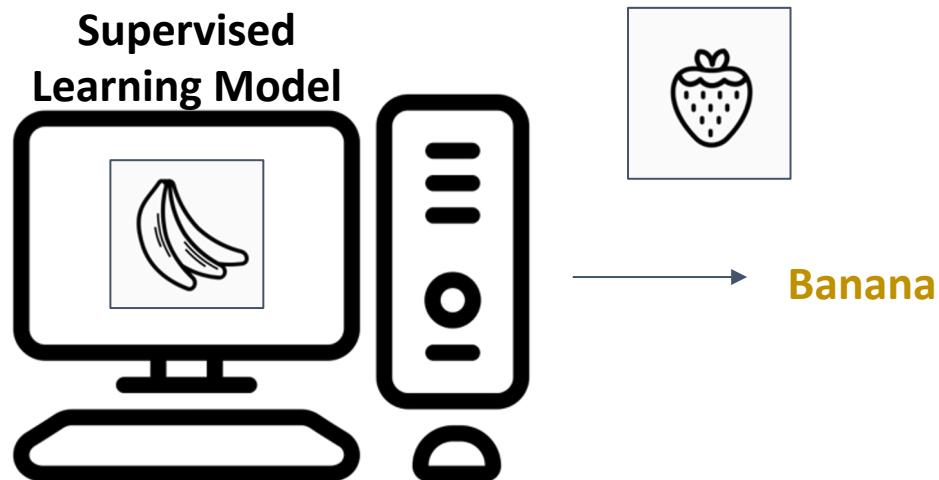
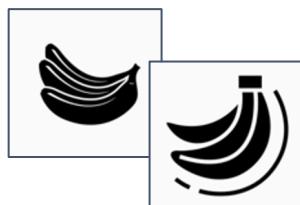
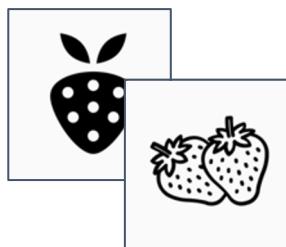
Supervised learning

- Have a bunch of labelled data,
want to label new data



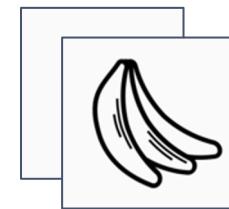
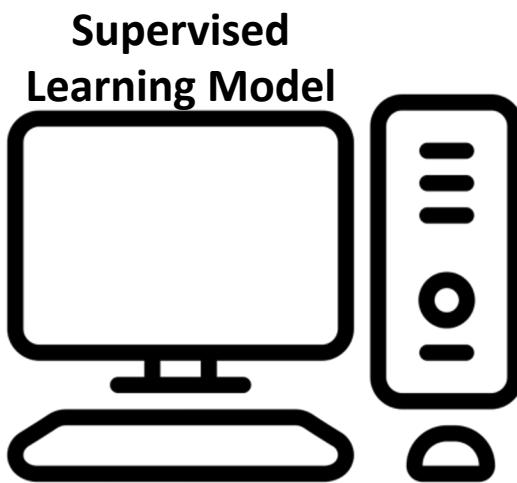
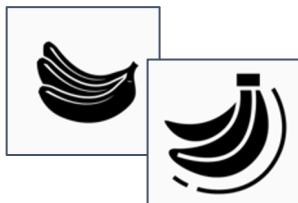
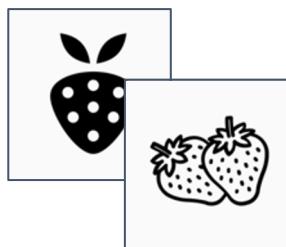
Supervised learning

- Have a bunch of labelled data,
want to label new data



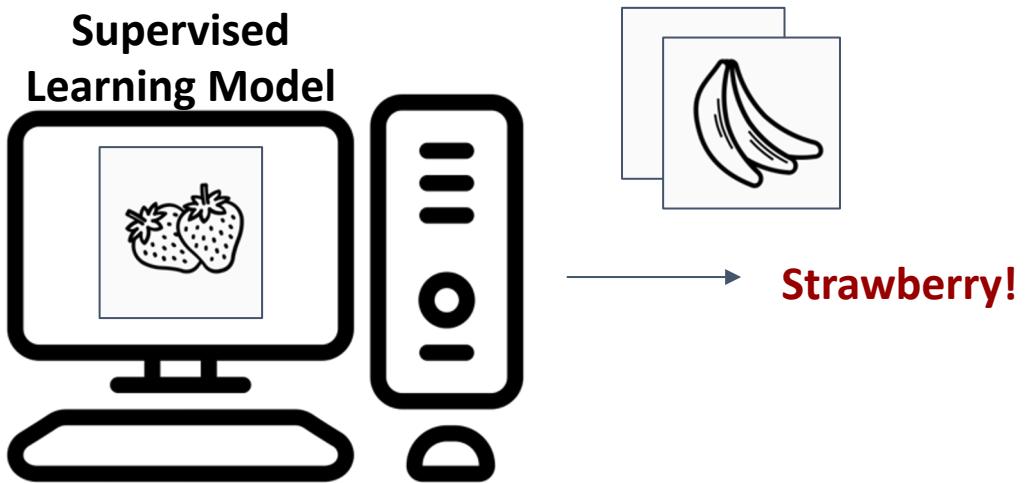
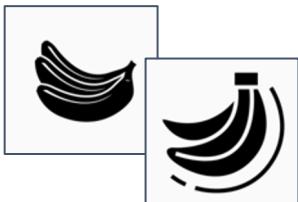
Supervised learning

- Have a bunch of labelled data,
want to label new data



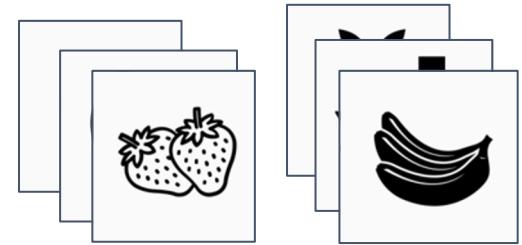
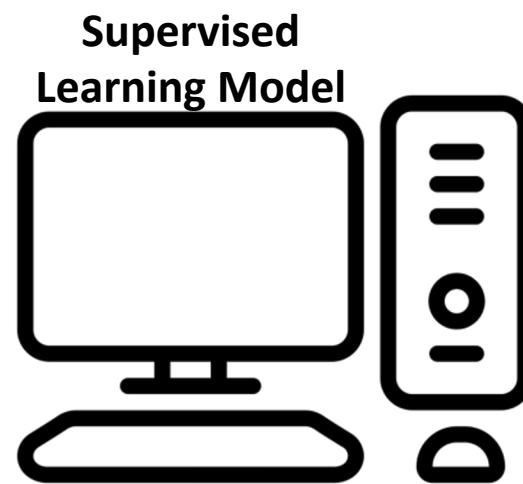
Supervised learning

- Have a bunch of labelled data,
want to label new data



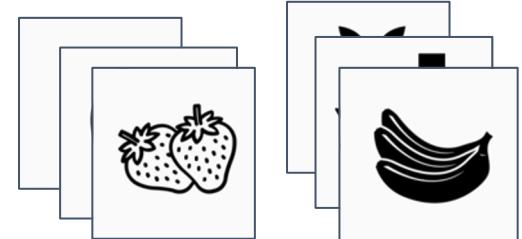
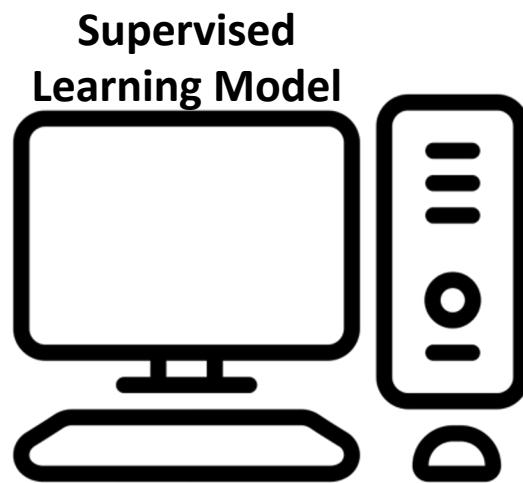
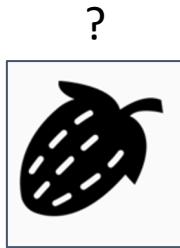
Supervised learning

- Have a bunch of labelled data,
want to label new data



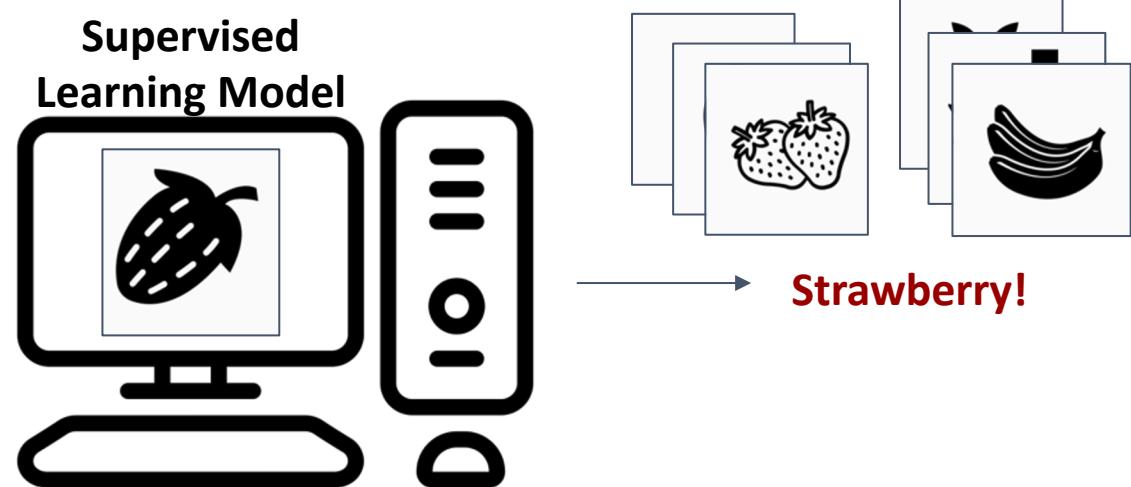
Supervised learning

- Have a bunch of labelled data,
want to label new data



Supervised learning

- Have a bunch of labelled data,
want to label new data





- Dataset contains features and a *label* or target
 - E.g. classify iris plants into the three different species
 - Mathematically, supervised learning observes examples of a random vector x and an associated label y
 - Goal is to learn to predict y from x
"supervised"
 - Term comes from view of the label y coming from an instructor or teacher

0000000000
1111111111
2222222222
3333333333
4444444444
5555555555
6666666666
7777777777
8888888888
9999999999



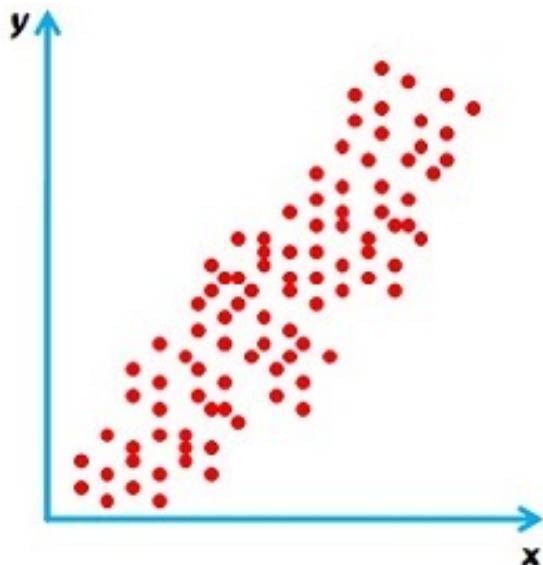
Supervised learning

- Input data space \mathcal{X}
- Output (label) space \mathcal{Y}
- Unknown function $f: \mathcal{X} \rightarrow \mathcal{Y}$
(mapping)
- Dataset $\mathcal{D} = \{(x_1, y_1), \dots (x_n, y_n)\}$ with $x_i \in \mathcal{X}$ and $y_i \in \mathcal{Y}$
- Finite $\mathcal{Y} \Rightarrow$ classification
- Continuous $\mathcal{Y} \Rightarrow$ regression



Regression

- We have a set of n observations (\mathbf{x}_i, y_i) with $y_i \in \mathbb{R}$
- Goal is to learn a function that predicts label y from sample \mathbf{x}

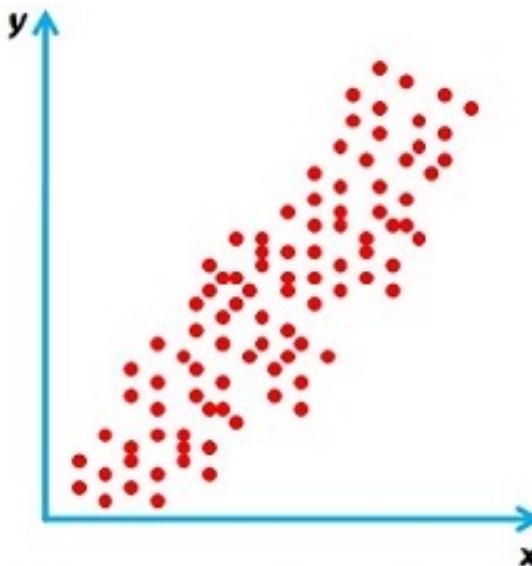


Surbhi S, 2016



Regression

1. Choose a ***model class*** of functions
model class of deep NN = last slide of 1st lecture (model selection)
2. Choose a performance measure to guide the selection of a function from the model class



Surbhi S, 2016

- We'll begin with a simple model class: linear functions



Linear Regression

- We want to fit a linear function to dataset $\mathcal{D} = \{(x_1, y_1), \dots (x_n, y_n)\}$
- $\hat{y} = \mathbf{w} \cdot \mathbf{x} = \mathbf{w}^T \mathbf{x}$
 - $\mathbf{w} \in \mathbb{R}^d$ is a vector of parameters
 - Line if $d = 1$
 - Plane if $d = 2$
 - Hyperplane in general d
- How do we determine \mathbf{w} ?
- Based on the performance measure P



Loss functions

- Labels are in \mathcal{Y}
 - Discrete (classification)
 - Continuous (regression)
- ***Loss function*** $L: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$
- L maps decisions/predictions to a cost.
 - $L(\hat{y}, y)$ is the penalty for predicting \hat{y} when the true label is y
- Standard choice for regression is ***mean squared error*** (MSE)
 - $L(\hat{y}, y) = (\hat{y} - y)^2$
- Absolute loss: $L(\hat{y}, y) = |\hat{y} - y|$



Empirical loss

实验 loss

- Parametric function $f(\mathbf{x}; \mathbf{w})$
 - E.g. $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$ for regression
- The *empirical loss* on a dataset \mathcal{D} is

$$L(\mathbf{w}; \mathcal{D}) = \frac{1}{n} \sum_i L(f(\mathbf{x}_i; \mathbf{w}), y_i)$$

- Least squares minimizes the empirical MSE
- We care about predicting labels for *new* samples
- Does empirical loss minimization help with that?

case of not help : new samples have different distribution with training samples.



Empirical vs. Expected loss

- Assumption: sample and label pairs (x, y) are drawn from the **unknown joint distribution** $p(x, y)$
- Data are sampled i.i.d.
- **Empirical loss** measured on training data

$$L(\mathbf{w}; \mathcal{D}) = \frac{1}{n} \sum_i L(f(\mathbf{x}_i; \mathbf{w}), y_i)$$

- Goal: minimize the **expected loss** (aka **risk**)

$$R(\mathbf{w}) = \mathbb{E}[L(f(\mathbf{x}; \mathbf{w}), y)]$$



Empirical risk (loss) minimization

- Empirical loss

$$L(\mathbf{w}; \mathcal{D}) = \frac{1}{n} \sum_i L(f(\mathbf{x}_i; \mathbf{w}), y_i)$$

- Risk

$$R(\mathbf{w}) = \mathbb{E}[L(f(\mathbf{x}; \mathbf{w}), y)]$$

- If training set is representative of $p(\mathbf{x}, y)$, then empirical loss is a suitable proxy for the risk
- Essentially, we estimate $p(\mathbf{x}, y)$ and $R(\mathbf{w})$ from the empirical distribution of the data

generalizing out of sample is a problem in NN
eg. can't recognize a cat with noise as cat



1. Choose a model class \mathcal{H} of functions $f: \mathcal{X} \rightarrow \mathcal{Y}$
 - Regression example: linear functions parameterized by \mathbf{w} :
$$f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x}$$
2. Select the function in \mathcal{H} that minimizes the empirical risk
 - Linear regression example: minimize empirical squared loss:
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i - y_i)^2$$
 - How do we find \mathbf{w}^* ?
 - Use calculus



Linear Regression Optimization

- Add an offset w_0 : $f(\mathbf{x}; \mathbf{w}) = \mathbf{w}^T \mathbf{x} + w_0$

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}_i + w_0 - y_i)^2$$

$$= \arg \min_{\mathbf{w}} L(\mathbf{w}; \mathcal{D})$$

- Set $\frac{\partial L(\mathbf{w}; \mathcal{D})}{\partial w_i} = 0$ for each i
- Results in $d + 1$ equations and $d + 1$ unknowns

$\mathbf{x} \in \mathbb{R}^d$



Linear Regression

- Switching to vector/matrix notation
 - X = data matrix with an additional column of 1's (for the offset)
 - \mathbf{y} = vector of labels
 - Prediction: $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$
- Resulting solution is
$$\mathbf{w}^* = (X^T X)^{-1} X^T \mathbf{y}$$
- Closed form solution!
- $(X^T X)^{-1} X^T$ is the Moore-Penrose pseudoinverse of X
- No closed form solution is available in most machine learning

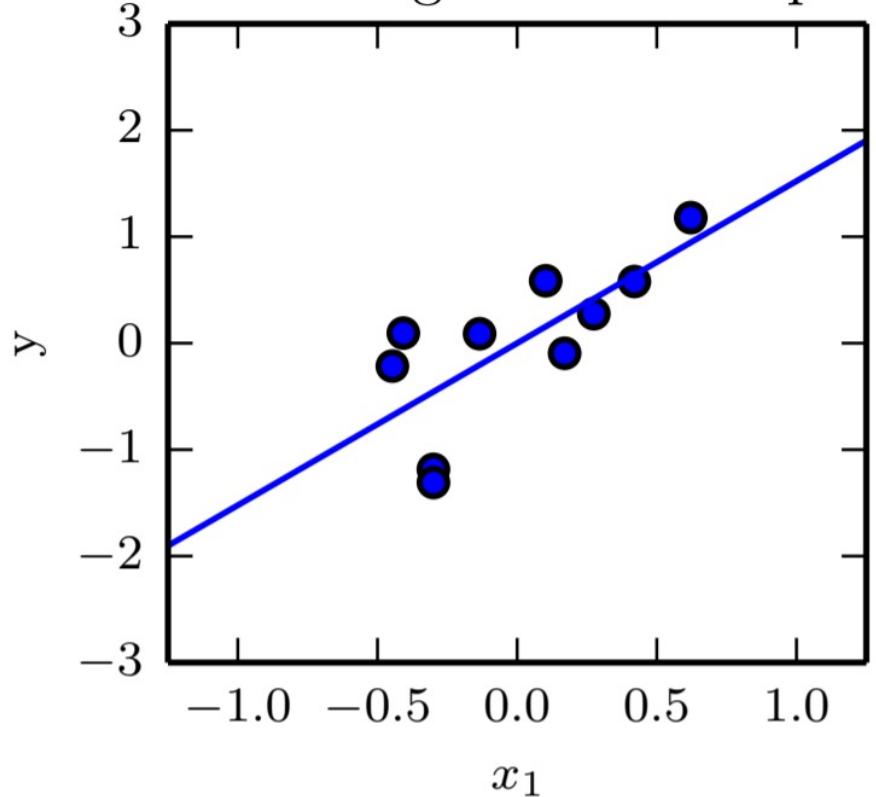
trade-off closed form solution vs function expressive
ness



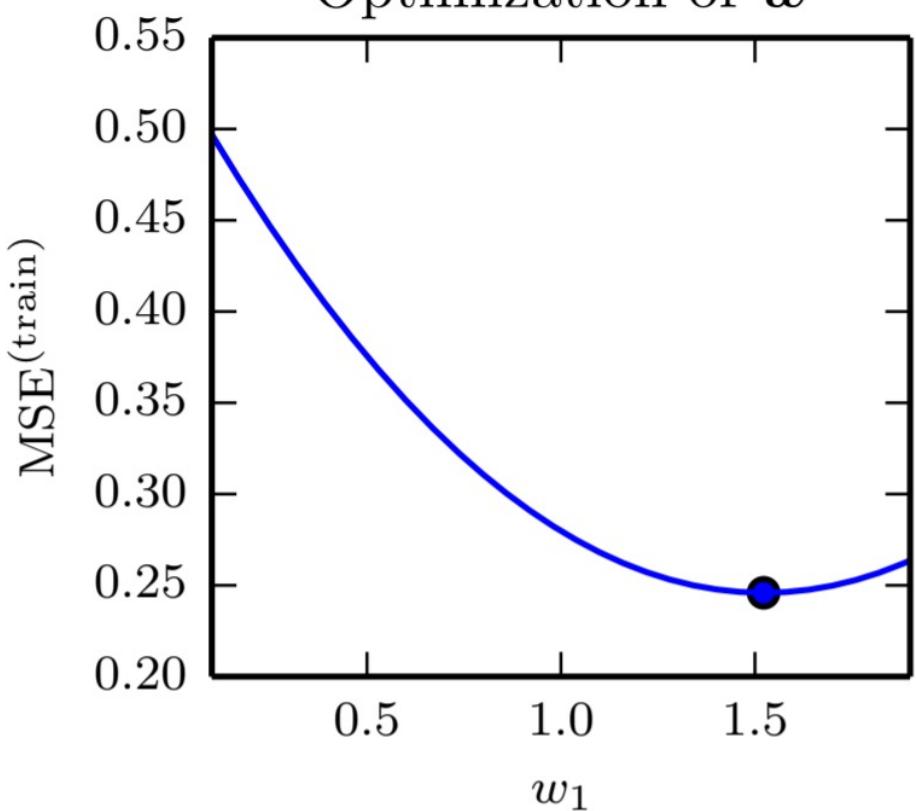
goal: find intercept and slope

✓ theory
optimization is in heart of ML

Linear regression example



Optimization of w



generalization



- Challenge of machine learning: algorithm must perform well on previously unseen inputs
 - Not just training inputs
- In other words, we want the algorithm to **generalize** well to unseen data
 - **generalization error** or the **test error measures error on unseen data**
- We also need to estimate the generalization error
 - Done typically by dividing the total amount of data into two sets: **training data** and **test data**
 - Model is trained on the training data and then tested on the test data
- Expected test error \geq expected training error

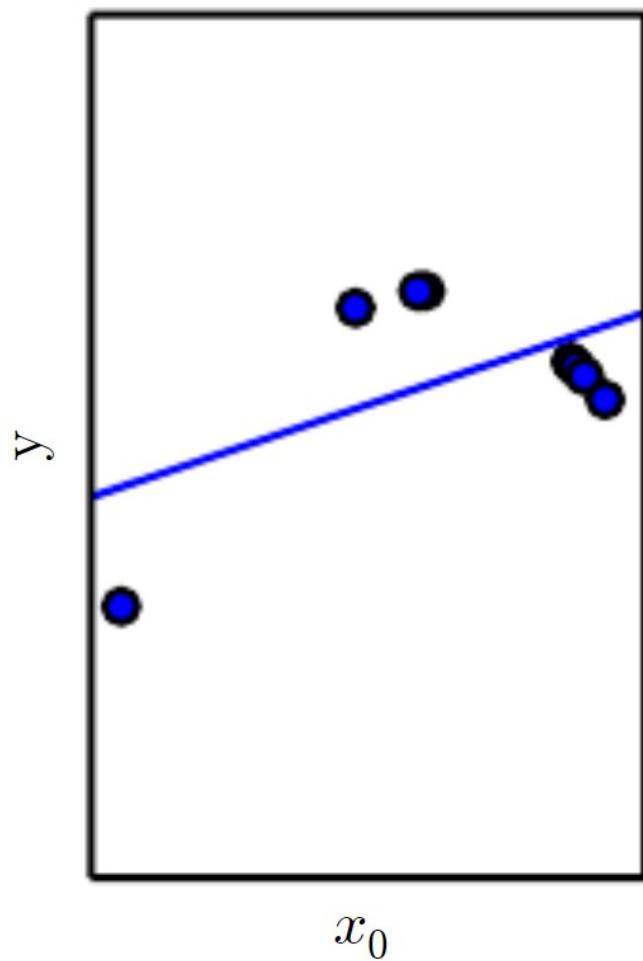
How much ?



Underfitting

Not choose right model class

- Many processes are not well modeled by linear functions
 - Biological systems
 - Stock market prices
 - Periodic signals
- In these cases, linear functions can ***underfit*** the data
 - Training error is very high
 - Test error will also be high





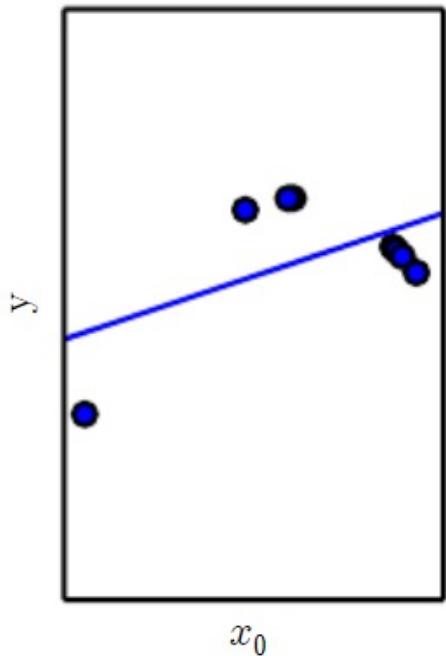
Model Complexity/Capacity

- Traditional Solution: increase **complexity** of the model
 - Model complexity is the # of independent parameters in the function model (i.e. “degrees of freedom”)
- We can increase the complexity of our linear function by adding polynomial terms
 - Example: $d = 1$
$$f(x; \mathbf{w}) = \sum_{j=0}^m w_j x^j$$
 - Still linear in \mathbf{w} \Rightarrow we can use linear regression!
 - Neural network solution may be to prefer certain types of solutions via regularizations

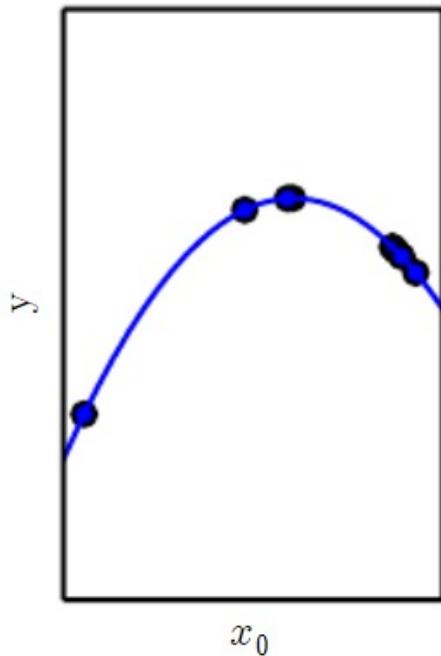


Nonlinear Regression

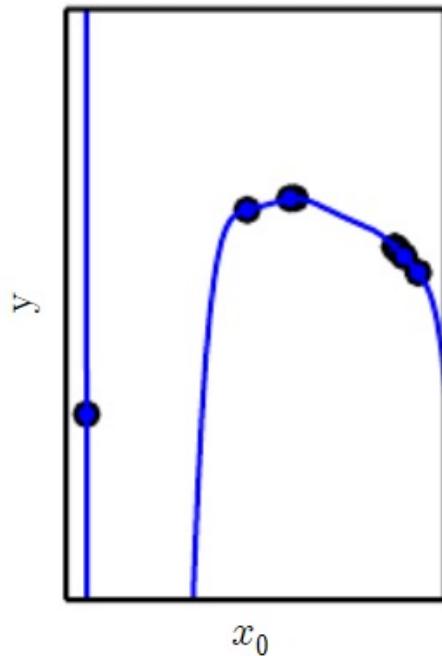
Underfitting



Appropriate capacity



Overfitting





Overfitting

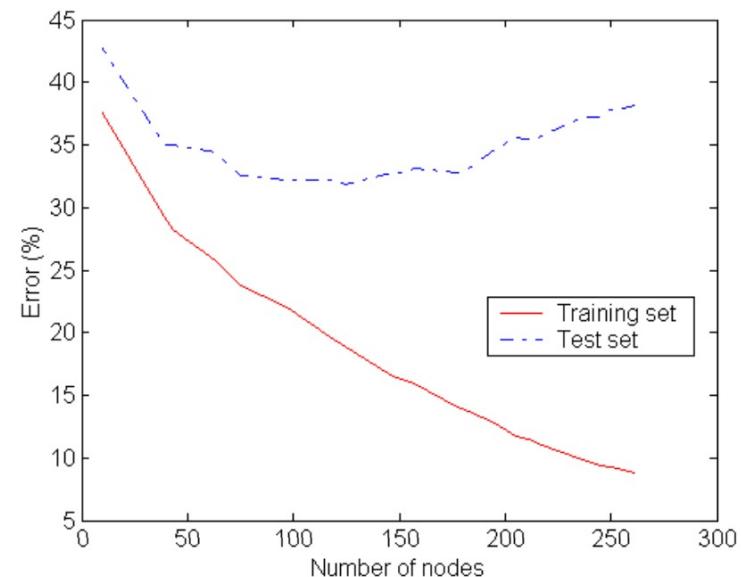
- Overfitting occurs when the gap between training error and test error is too large

- The model is capturing unique characteristics of the training set
- The model does not generalize well to new samples
- Occurs when the complexity is too large

- Avoiding overfitting

- Generally requires the use of another test set for evaluation
- Traditionally: decrease complexity

another solution: regularization



Inductive bias for simplicity



0 coefficients

a preference for more

one kind of regularization

- Include a penalty on model complexity directly in the cost function:

$$\sum_i L(f(\mathbf{x}_i; \mathbf{w}), y_i) + \text{\#parameters}$$

- The addition of the penalty is known as **regularization**
- Regularization: “any modification we make to a learning algorithm that is intended to reduce its generalization error but not its training error.” (Goodfellow et al., 2016)
 - Lots of different penalties can be included in regularization
 - Regularization is useful/important in deep learning (and machine learning in general)
- Bias does not have to be directly complexity

regularization
express a preference for simplicity in solution , added
simplicity { fewer params
smaller range of param values
inductive bias



Regularization

- Ridge regression: penalize with L2 norm

$$\mathbf{w}^* = \arg \min \sum_i L(f(\mathbf{x}_i; \mathbf{w}), y_i) + \lambda \sum_{j=1}^m w_j^2$$

prefer smaller coefficient

- Closed form solution exists $\mathbf{w}^* = (\lambda I + X^T X)^{-1} X^T \mathbf{y}$ *for linear regression*
- LASSO regression: penalize with L1 norm

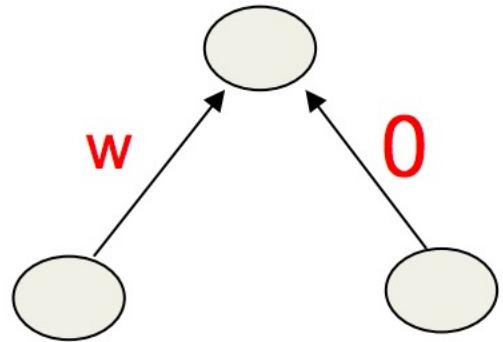
$$\mathbf{w}^* = \arg \min \sum_i L(f(\mathbf{x}_i; \mathbf{w}), y_i) + \lambda \sum_{j=1}^m |w_j|$$

- No closed form solution but still convex (optimal solution can be found)
by gradient descent

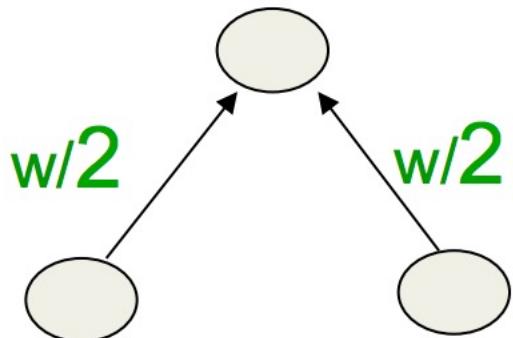
prefer 0 coefficient



L1 vs L2



- Prefers to share smaller weights
- Makes model smoother
- More Convex





L1 vs L2

L2 regularization	L1 regularization
Computational efficient due to having analytical solutions	Computational inefficient on non-sparse cases
Non-sparse outputs	Sparse outputs
No feature selection	Built-in feature selection

smaller coefficient made function less fluctuate and smoother

Other regularization

{ increase size of training data
early stopping ~~is~~ good for GAN
drop out



L_p Regularization

$$L_p = \left(\sum_i w_i^p \right)^{1/p}$$

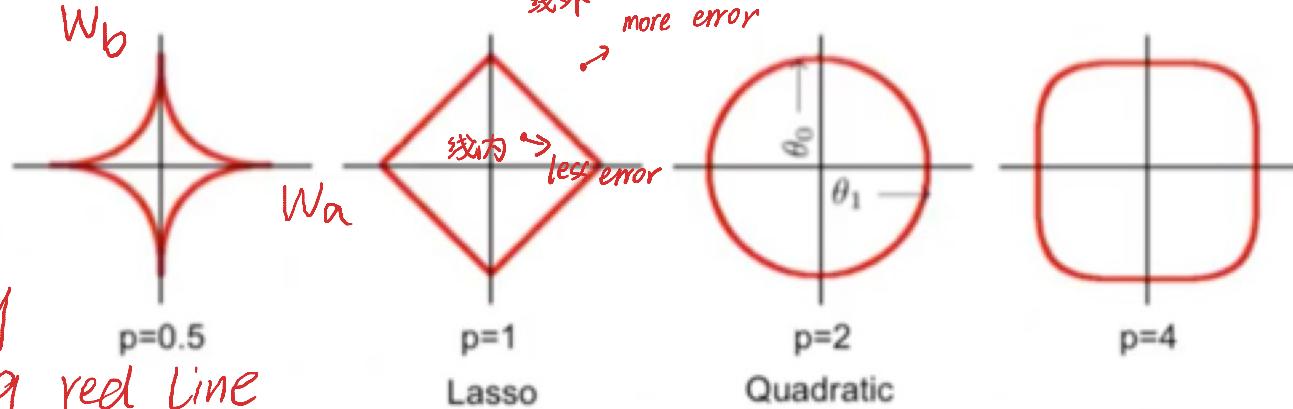
Isosurfaces

等值面

same penalty

value along red line

assume 2 params a, b



only care about which param has max value

$P \rightarrow \infty$

$$L_p = \max(w_i) \quad (i=1, \dots, d)$$

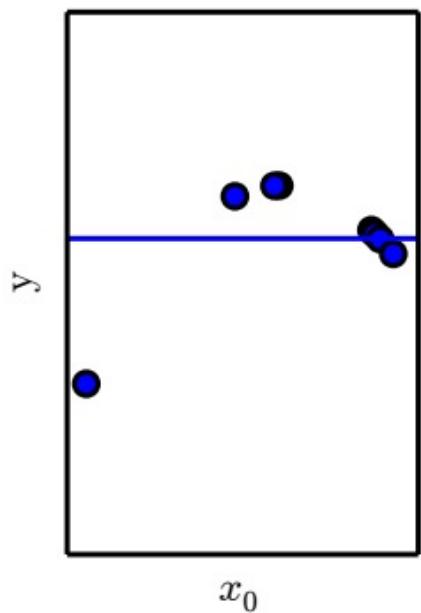
goal: # param ↓ As $p \rightarrow 0$ you get a counting penalty

$$L_p = \|\sum w_i\|_0 = \# w$$

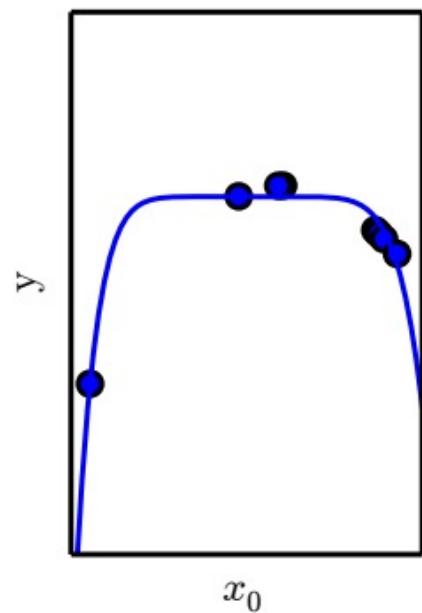
Goal: so try to bring down weight have max value



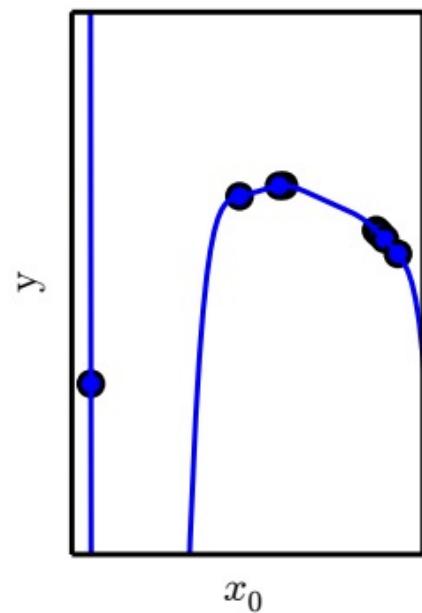
Underfitting
(Excessive λ)



Appropriate weight decay
(Medium λ)

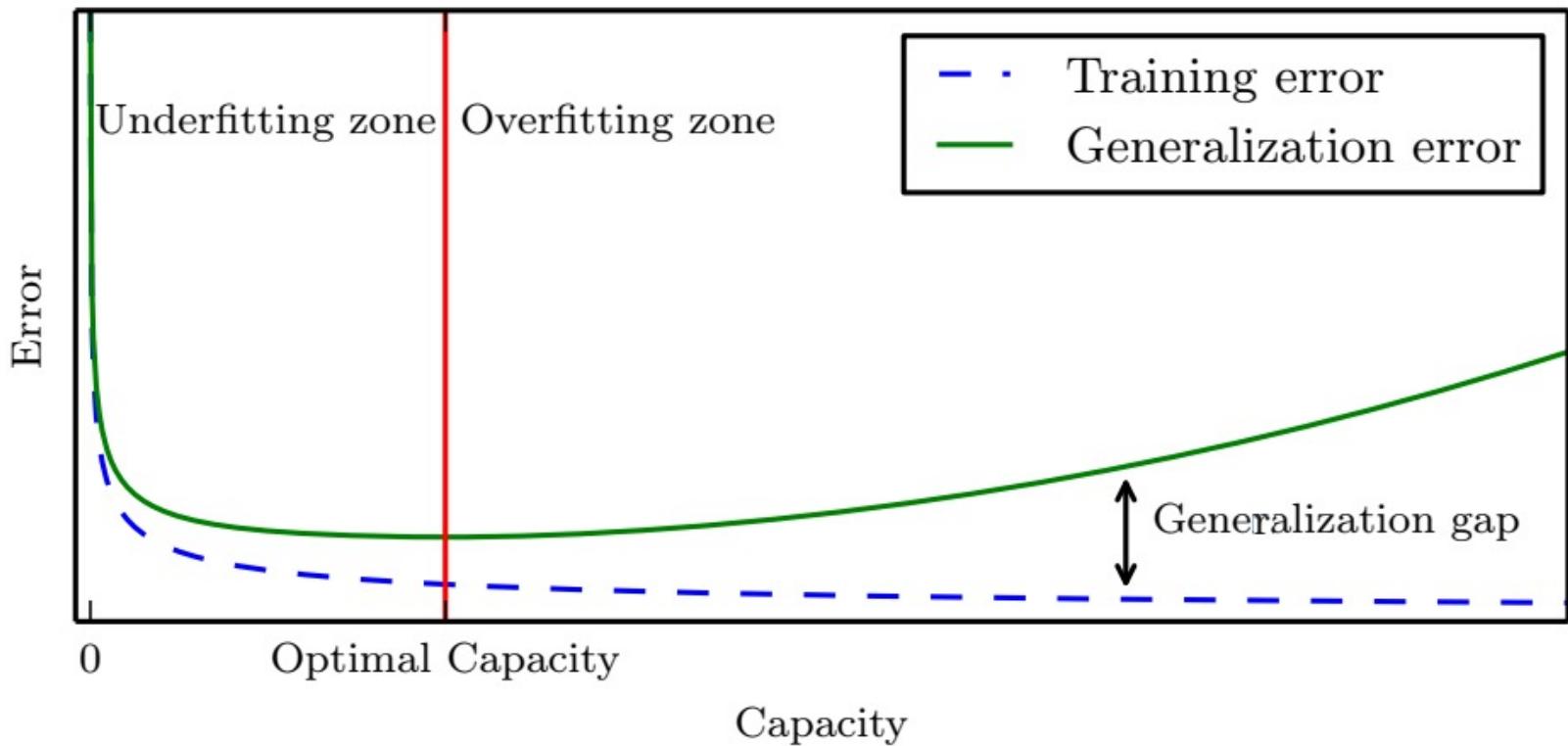


Overfitting
($\lambda \rightarrow 0$)





Underfit/Overfit Paradigm

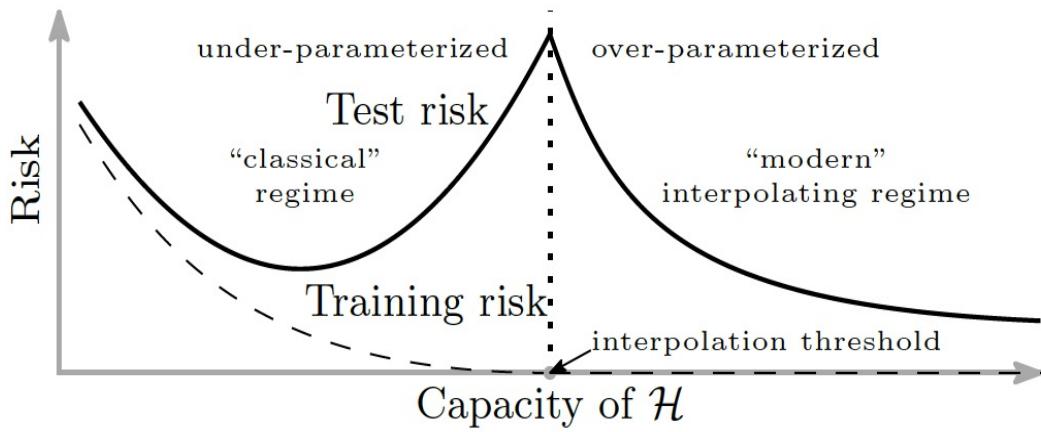
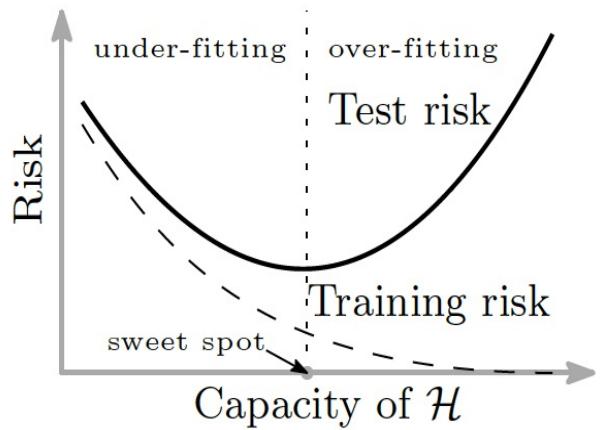


Conventional Wisdom Doesn't work!





New “Interpolation” Regime





k -nearest neighbor classifier

an example of non-NN over-parameterized model

here increase #params is increase k

- No training required
- 1. Find the k -nearest neighbors of the test point x in the training data
 - Call this set $\mathcal{N}_k(x)$
- 2. Choose \hat{y} according to a majority vote of the corresponding labels of the neighbors in $\mathcal{N}_k(x)$

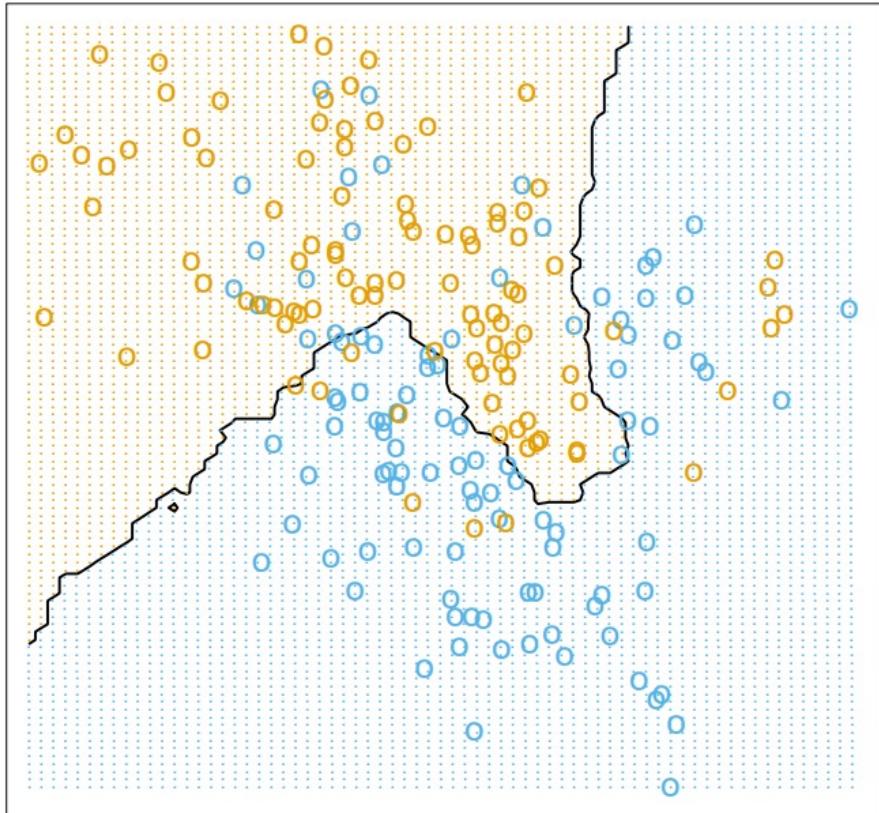


Interpolation Regime

插值机制

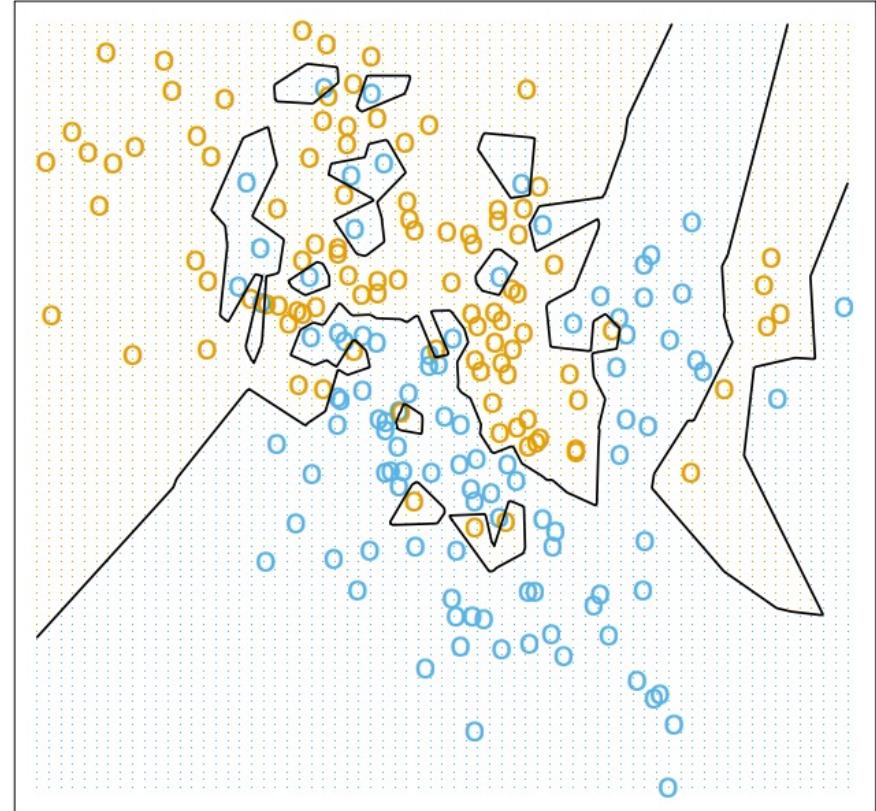
smoother

15-Nearest Neighbor Classifier



overfitting

1-Nearest Neighbor Classifier



k -nearest neighbor classifier

Hastie et al., 2009



Measuring test error

- Different from training error
- How do we select the data for estimating the two?
- Common approaches:
 - **Holdout validation:** Randomly split the data into two disjoint sets. Train the model on one set and compute test error on the other.
 - **Random subsampling:** Repeat the holdout validation k times independently and compute the mean test error
 - Typically these approaches are not very robust

Caveat on NN

first train { training set: A
 | test set: B \Rightarrow 2 different NNs with different weights and biases

second train { training set: B hard to report p-value
 | test set: A



Cross Validation

- Cross validation averages evaluation scores over several different splits to provide a more realistic assessment
- Types of cross validation
 - **Leave-one-out:** Use all data points except one for training and validate on the remaining point.
 - **2-fold cross validation:** Randomly split the data into two halves and run two validation iterations, each using one half for training and the other for validation
 - **k-fold cross validation:** Randomly partition the data into k equal size parts and run k validation iterations, each using one part for validation and the rest for training
- In all cases, the validation scores are averaged
- Each data point is only used once in testing but many times in training

another layer, randomization of training data



Further reading

- Chapter 5 in Goodfellow book
- *Elements of Statistical Learning* by Hastie, Tibshirani, & Friedman
- *Pattern Recognition and Machine Learning* by Bishop

