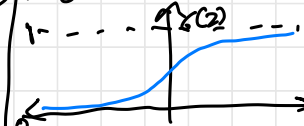
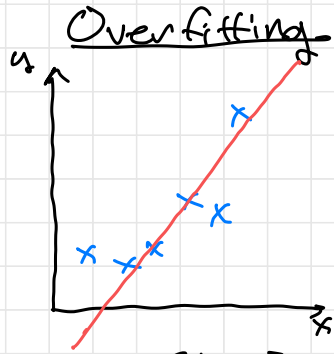


1/26/2023 : Overfitting, Regularization

Review of linear supervised learning methods so far
 learn from dataset of (x, y) pairs

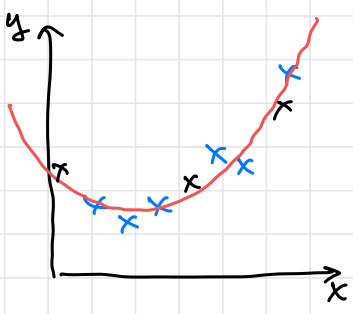
input \nearrow correct answer "supervision" \nwarrow

	Linear Regression	Logistic Regression	Softmax Regression
Task	Regression $y \in \mathbb{R}$	Binary classification $y \in \{-1, 1\}$	multiclass classification $y \in \{1, 2, 3, \dots, C\}$
Parameters (what to learn)	$w \in \mathbb{R}^d$ dimension of x	$w \in \mathbb{R}^d$	$w^{(1)}, \dots, w^{(C)} \in \mathbb{R}^d$ total $C \cdot d$ params
Probabilistic Story	$y \sim \text{Normal}(w^T x, \sigma^2)$ mean \nearrow	$p(y=1 x) = \sigma(w^T x)$ 	$p(y=j) = \frac{\exp(w^{(j)T} x)}{\sum_{k=1}^C \exp(w^{(k)T} x)}$ Normalizes to probability distribution
How to get loss function measures how bad any choice of parameters is	Maximum Likelihood Estimation (MLE) maximize probability of data = $\prod_{i=1}^n p(y^{(i)} x^{(i)}; w)$ with respect to w \Leftrightarrow minimize negative log likelihood = $-\sum_{i=1}^n \log p(y^{(i)} x^{(i)}; w)$		
How to minimize loss	Gradient descent OR Normal Equations	Gradient descent 1 st -order OR Newton-Raphson Method 2 nd order (also uses 2 nd derivs)	



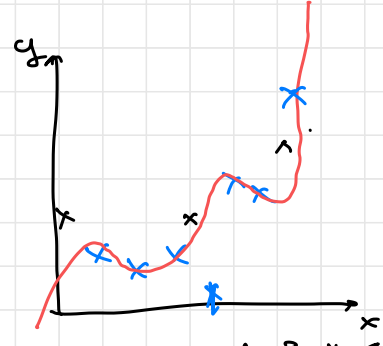
Features = $[1, x]$

Too simple
"underfitting"



Features = $[1, x, x^2]$

★ Best fit



Features = $[1, x, x^2, x^3, x^4, x^5]$

Sensitive to fluctuations
in training
"overfitting"

Not THE SAME

Lowest training loss

Dataset Splits: ALWAYS have 3 datasets with disjoint examples

Training set

Use it to choose parameters
(e.g. run linear regression with this dataset)

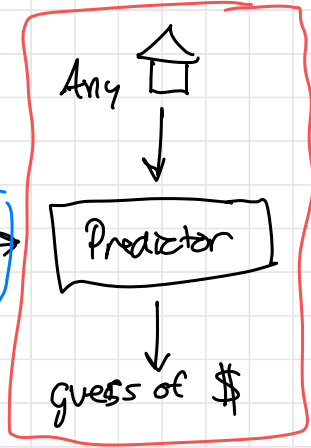
Development Set
("Validation")

Test set

Evaluate how well our model does

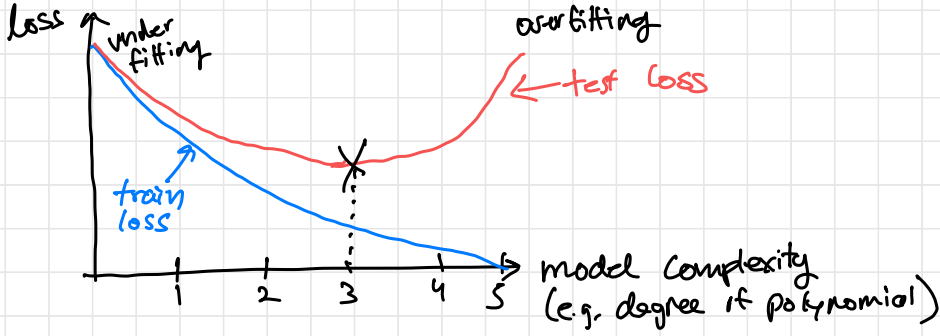


Just to build the predictor



What matters:
accuracy of predictions on unseen examples

Test set simulates model behavior on unseen examples



What is development set for?

A: To choose hyperparameters

Any setting of learning algorithm

- Learning rate
- When to stop training
- Which features?

Rule: Choose hyperparameters to minimize development set loss,

(In contrast, parameter is chosen by the learning algorithm)

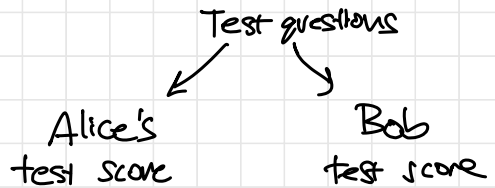
only evaluate on test set at very end.

degree	<u>dev</u> loss	final evaluation <u>test</u> loss
1	100	100
2	51	50
3	50	50
4	49	50
5	75	75

- ① Train 5 models
- ② Evaluate each on dev set
- ③ Pick the model that's best on dev
- ④ Evaluate that model on test set

Announcements

- HW 1 released Feb 7
- HW 0 grades out



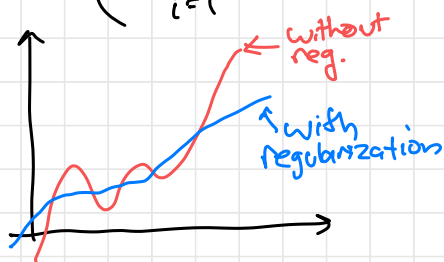
"Correlation does not imply Causation"

Regularization: A way to reduce overfitting by preferring "simpler" functions

L2 Regularization: Encourage ^{square of} L2 norm of parameters to be small $\sum_{j=1}^d w_j^2 = \|w\|^2$

e.g. linear regression

$$L(w) = \left(\frac{1}{n} \sum_{i=1}^n (w^T x^{(i)} - y^{(i)})^2 \right) + \lambda \|w\|^2$$



Some constant hyperparameter
 $\lambda = 0 \Rightarrow$ no regularization
 λ large \Rightarrow strong regularization

How does this change the gradient

gradient of $\lambda \|w\|^2 = 2 \lambda w$

during G.D., you subtract $\eta \cdot 2 \lambda w$

learning rate

"weight decay"

L1 Regularization: add $\lambda \|w\|_1$ to objective

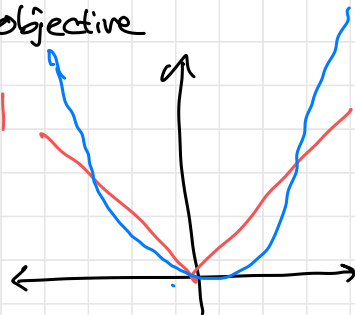
$$= \sum_{j=1}^d |w_j|$$

Gradient for L1:

$$\frac{d}{dw_j} \lambda \|w\|_1 = \lambda \text{sgn}(w_j)$$

So full gradient is $\lambda \begin{bmatrix} \text{sgn}(w_1) \\ \vdots \\ \text{sgn}(w_d) \end{bmatrix}$

$$\text{sgn}(z) = \begin{cases} 1 & \text{if } z > 0 \\ -1 & \text{if } z < 0 \\ 0 & \text{if } z = 0 \end{cases}$$



constant sized step towards 0

vs $2 \lambda w$ from L2

if w close to 0, you take very small step

L_1 regularization has a sparsifying effect
(leads to sparse w)

means many entries = 0

L_2 avoids very big entries of w