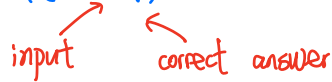


Review:

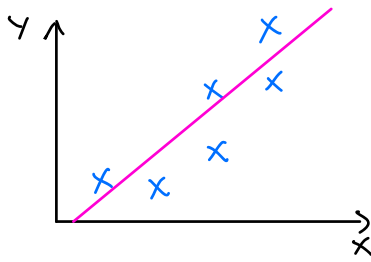
linear supervised learning

learning from data  $(x, y)$



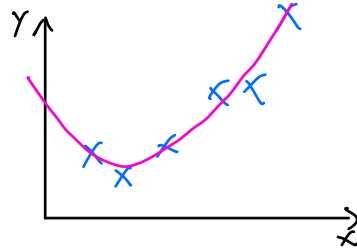
	linear regression	logistic regression	softmax regression
Task	regression $y \in \mathbb{R}$	binary classification $y \in \{-1, 1\}$	multi-class classification $y \in \{1, 2, \dots, C\}$
Parameters	$w \in \mathbb{R}^d$ $x \in \mathbb{R}^d$	$w \in \mathbb{R}^d$	$w^{(1)}, \dots, w^{(C)} \in \mathbb{R}^d$ c-d params
Probabilistic story	$y \sim N(w^T x, \sigma^2)$ mean	$p(y=1 x=x) = \sigma(w^T x)$ ... 	$p(y=j) = \frac{\exp(w^{(j)T} x)}{\sum_{k=1}^C \exp(w^{(k)T} x)}$ normalises to prob. dist.
How to get loss functions	MLE, maximize prob. of data = negative log likelihood $\prod_{i=1}^n p(y^{(i)}   x^{(i)}; w)$ w.r.t. $w \Leftrightarrow$ minimize NLL		
How to minimize losses	gradient descent or normal eqns	gradient descent 1st-order or Newton-Raphson method 2nd-order	

## Overfitting



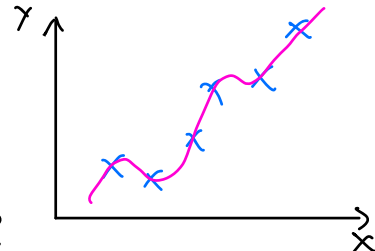
features =  $[1, x]$

too simple  
"underfitting"



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

★ best fit



features =  $[1, x, x^2, x^3, x^4]$

sensitive to fluctuations in training  
"overfitting"

## Data Splits

Always have 3 disjoint datasets

### Training

Use it to choose parameters

### Validation

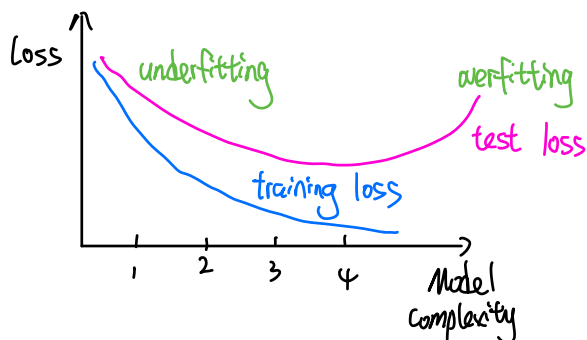
Choose hyperparameter  
any setting of learning algorithm

- learning rate
- features
- when to stop training

### Test

Evaluate how well our model performs

parameters chosen by the learning algorithm



Rule: choose hyperparameters to minimize validation loss, only evaluate on test set at very end.

degree	(validation)	
	development loss	test loss
1	100	100
2	51	50
3	50	50
4	49	50
5	75	75

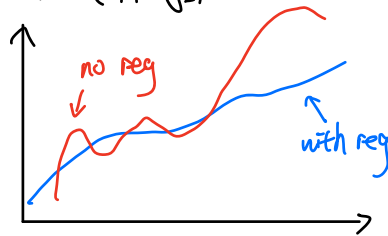
- ① train 5 models
- ② evaluate each on dev
- ③ pick the best based on ②
- ④ evaluate that model on test set

**Regularization** (to reduce overfitting prefer "simpler" functions)

**L2 Regularisation**  $\leftarrow \sum_{j=1}^d w_j^2 = \|w\|^2$   
 encourage square of L2 norm to be small

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 hyperparameters  
 $\lambda = 0 \Rightarrow$  no reg  
 $\lambda$  large  $\Rightarrow$  strong reg

How does this change the gradient

$$\text{gradient of } \lambda \|w\|^2 = 2\lambda \|w\|$$

doing GD, subtract  $\eta \cdot 2\lambda w$

$\uparrow$   
learning rate

**L1 Regularization**

add  $\lambda \|w\|$  to objective  $\|w\| = \sum_{j=1}^d |w_j|$

$$\text{gradient of } L_1: \frac{\partial}{\partial w_j} \lambda \|w\| = \lambda \text{sgn}(w_j)$$

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

constant size step towards 0  $\rightarrow$

VS  $2\lambda w \leftarrow$  if  $w=0$ , take very small step

$L_1$  reg has a sparsifying effect  
(leads to sparse  $w$ )  
many entries = 0

$L_2$  reg avoids very big entries of  $w$