

## Naive Bayes

$$P(Y=1 | x_1=c, x_2=0) = \frac{P(x_1=c | Y=1) P(x_2=0 | Y=1) P(Y=1)}{P(x_1=c, x_2=0 | Y=1) + P(x_1=c, x_2=0 | Y=0)}$$

## inductive bias

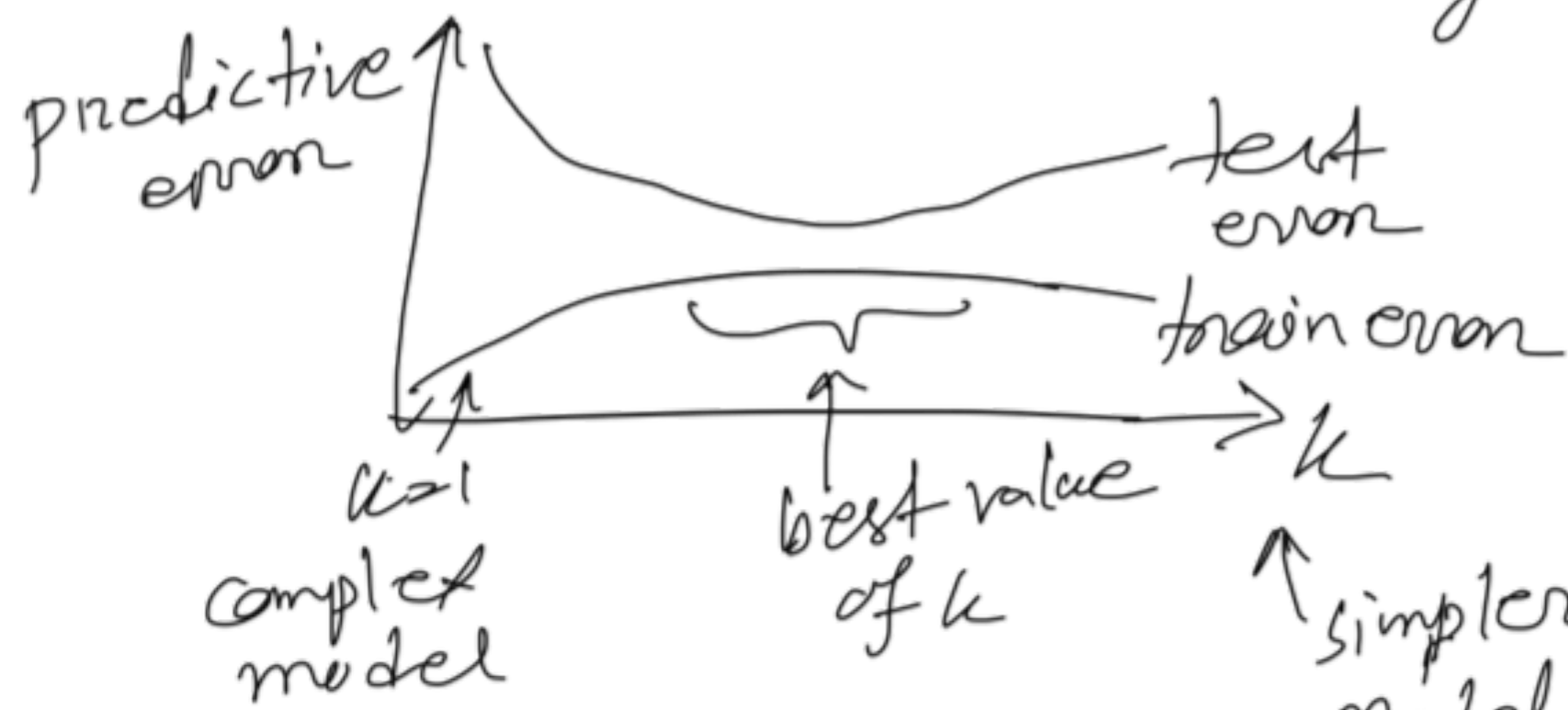
set of assumptions on the learner before it sees any data

## Curse of dimensionality in KNN

In higher dimension, distance from farthest and nearest points to the query point is same.

## KNN

- as  $k \downarrow \rightarrow$  overfit;  $k \uparrow \rightarrow$  smoother boundaries
- regression  $\rightarrow$  avg over  $y$  values
- classification  $\rightarrow$  majority voting



- as  $m$  increases, the optimal value of  $k$  increases
- use validation error to select best  $k$ .

## Regression

$$J(\theta) = \frac{1}{2m} \sum_i (y_i - \theta x_i^T)^2$$

$$\nabla J(\theta) = -\frac{1}{m} \sum_i [y_i - \theta x_i^T] \cdot [x_i^0, x_i^1, x_i^2, \dots]$$

error magnitude      sensitivity to each param.

$$= -\frac{1}{m} \sum_i [Y^T - \theta X^T] \cdot [X]$$

- for min. MSE,  $\nabla J(\theta) = 0$
- $\Rightarrow Y^T X - \theta X^T X = 0$
- $\Rightarrow \theta = Y^T X (X^T X)^{-1}$

- issues with MSE choice: sensitivity to outliers
- L1 loss =  $\sum_i (y_i - \hat{y}_i)$  L2 loss =  $\sum_i (y_i - \hat{y}_i)^2$
- Best MSE constant predictor = mean of the training points
- Best MSE linear predictor = pass through the points (if 2) or three midpoints (if  $\geq 2$ )



- Bias, Variance both depend on Learner, training procedure, training size

- Detect overfitting/select model: use validation

- Cross Validation

adv: lets us use more validation data so, less noisy estimate of test data

disadv: more work, trains  $k$  model doesn't evaluate any particular predictor

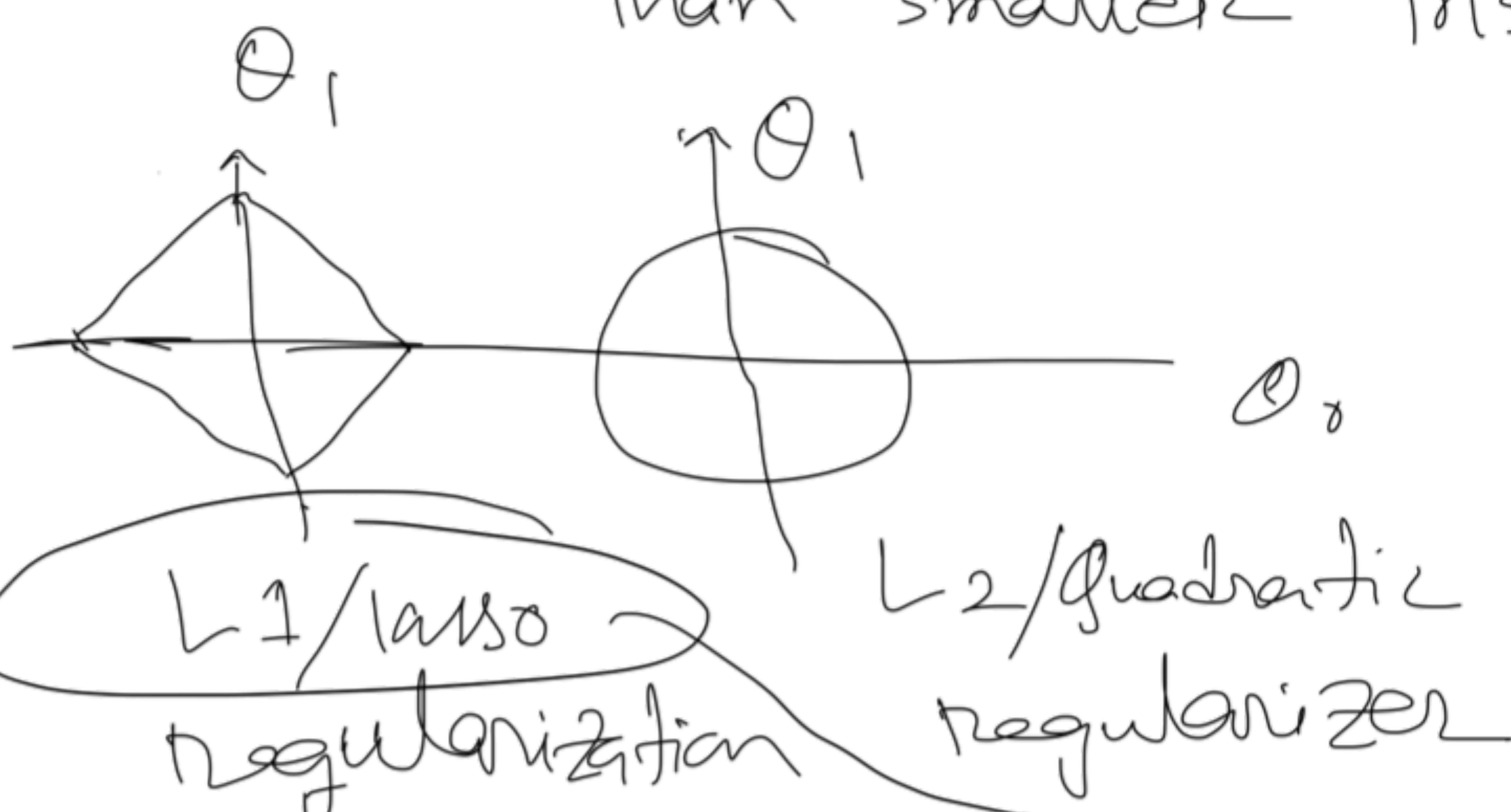
- keep degree fixed, increase data size  $\rightarrow$  low variance

- regularize when overfit/complexity high

$$J(\theta) = \frac{1}{2} (Y - \theta X^T) \cdot (Y - \theta X^T)^T + \alpha \theta \theta^T$$

regularizer

$\alpha$  large  $\Rightarrow$  smaller  $\theta$  more preferable than smaller MSE



Smaller $\alpha$	Larger $\alpha$
penalize small non-zero params more	less
a few big params	No big params
	prefers lots of small weights

generates sparser soln than L2.

Data increasingly separable in higher dimension

## Perceptron Algo

while - done  
for each data point  $i$ :

$$\hat{y}_i = \text{sign}(\theta \cdot x_i)$$

$$\theta = \theta - (y_i - \hat{y}_i) x_i$$

In logistic regression, we interpret  $\sigma(\theta x^T)$  as a prob. that  $y=1$ . The NLL loss

$$J(\theta) = \frac{1}{m} \sum_i (y_i \log(\sigma(\theta x_i^T)) + (1 - y_i) \log(1 - \sigma(\theta x_i^T)))$$

also, replace 0/1 loss with surrogate loss.

Multiclass perceptron:

while - done:

for each data point

$$f_i = \arg \max_c (\theta_c \cdot x_i)$$

$$\theta_f = \theta_f - \alpha x_i$$

$$\theta_y = \theta_y + \alpha x_i$$

[if prediction  $f$  matches true class  $y$ , then no update]

## NNet

activation:

$$\text{logistic} = \frac{1}{1 + \exp(-z)}$$

$$\text{ReLU} = \max(0, z)$$

## SVM

primal form:

$$w^* = \arg \min_w \sum_j w_j^2$$

$$\text{s.t. } y_i (w x_i^T + b) \geq +1$$

adv: when lots of data

Dual:

$$\max_{\alpha \geq 0} \sum_i \left( \alpha_i - \frac{1}{2} \sum_j \alpha_j y_j y_i (x_i \cdot x_j) \right)$$

$$\text{s.t. } \sum_i \alpha_i y_i = 0$$

adv: when lots of features efficient

## Perceptron

$m \uparrow \rightarrow$  bias  $\downarrow$  False  
 $n \uparrow \rightarrow$  variance  $\uparrow$  True

SGD adv: better for large  $m$  may avoid local optima

BRD adv: easier to set step size monotonic easier to gauge convergence  
 $\downarrow$   
may stuck in local optima

## On overfitting

$m \uparrow - \downarrow$  } lin. Regression  
 $n \uparrow - \uparrow$

$m \uparrow - \uparrow$  } KNN  
 $n \uparrow -$  No effect  
Rescale = u

# hidden nodes -  $\uparrow$  } Neural Net

# sigmoid  $\rightarrow$  Relu  
- No effect

kernel lin  $\rightarrow$  poly  
-  $\uparrow$  } SVM

grad descent/SG  
- No effect

$R \uparrow - \downarrow$  / No. of } soft SVM

lin Reg  $\Rightarrow$  Regularization  $\uparrow$   
 $\rightarrow$  variance  $\downarrow$

feat. transform  $\Rightarrow$  No effect on VC Dim

lin. classifier:  
overfitted,  $m \uparrow$   
 $\Rightarrow$  train err  $\uparrow$   
test err.  $\downarrow$   
underfitted,  $m \uparrow$   
 $\Rightarrow$  train err same  
test err. same

but increasing  
# of features  
increases VC dim.