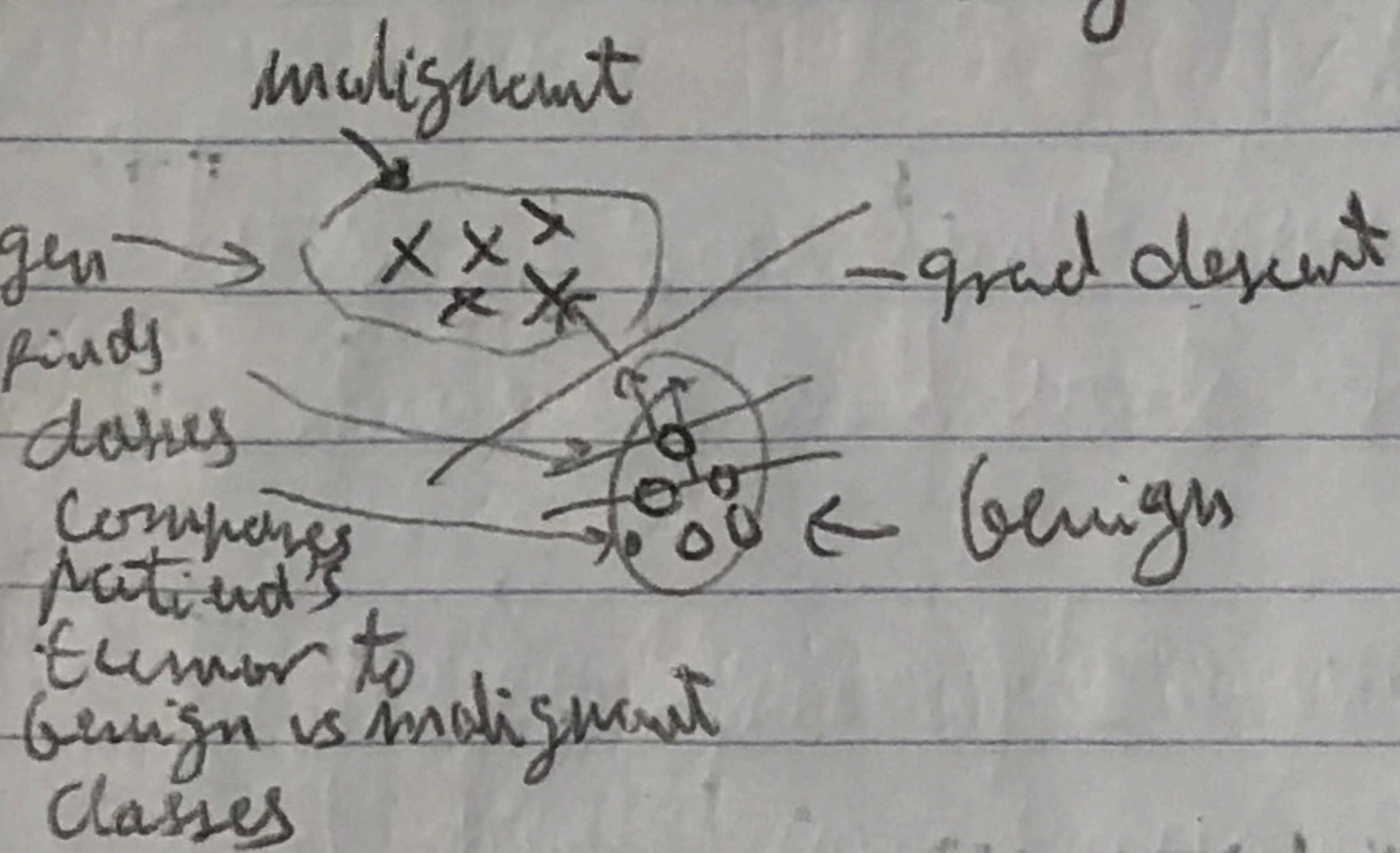


## Lecture 5

Today: Generative Reasoning Algorithms

- Gaussian Discriminant Analysis (GDA)
- Generative vs Discriminative Comparison
- Naive Bayes



Discriminative: Learn  $p(y|x)$  (or learn  $\hat{h}_\theta(x) = \{0, 1\}$  directly)

Generative: learns  $p(x|y)$   $p(y)$

$$\text{Bayes Rule: } p(y=1|x) = \frac{p(x|y=1)p(y=1)}{p(x)}, \quad p(x) = p(x|y=1)p(y=1) + p(x|y=0)p(y=0)$$

Gaussian discr. analysis (GDA)

Suppose  $x \in \mathbb{R}^n$  (drops  $x_0=1$  convention). Assume  $p(x|y)$  is Gaussian.

$$z \sim \mathcal{N}(\vec{\mu}, \Sigma), z \in \mathbb{R}^n, E(z) = \mu \text{ & } \text{Cov}(z) = E((z-\mu)(z-\mu)^T) = E_{zz^T} - E_z E_z^T$$

$$p(z) = \frac{1}{\sqrt{2\pi} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right]$$

GDA model; Assume:

$$p(x|y=0) = \frac{1}{\sqrt{2\pi} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu_0)^T \Sigma^{-1} (x-\mu_0)\right) \quad \text{Parameters: } \mu_0, \mu_1, \Sigma, Q.$$

$$p(x|y=1) = \frac{1}{\sqrt{2\pi} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x-\mu_1)^T \Sigma^{-1} (x-\mu_1)\right) \quad \text{Same, assumed typically}$$

$$P(y) = \phi^y (1-\phi)^{1-y} \quad P(y=1) = \phi.$$

Training set:  $\{(x^i, y^i)\}_{i=1}^m$

Joint likelihood for generative:

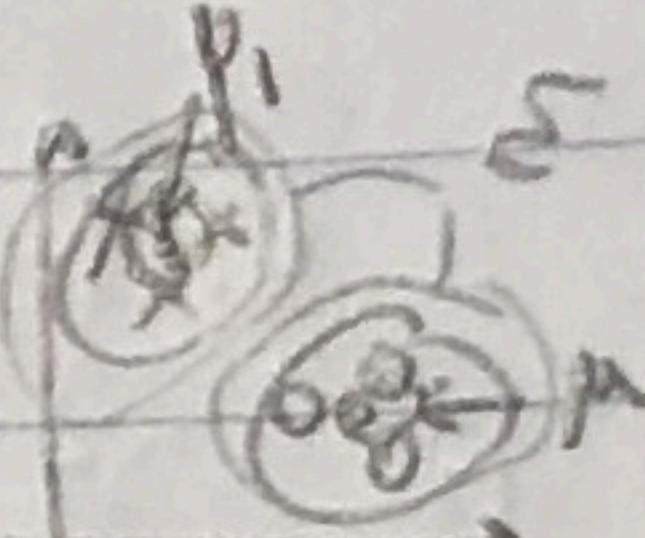
$$L(\phi, \mu_0, \mu_1, \Sigma) = \prod_{i=1}^m p(x^i, y^i; \phi, \mu_0, \mu_1, \Sigma) = \prod_{i=1}^m p(x^i | y^i) p(y^i)$$

& Discriminative:  $d(\theta) = \prod_{i=1}^m p(y^i | x^i; \theta)$  (conditional likelihood)

Maximum Likelihood Estimation:

$$\max_{\phi, \mu_0, \mu_1, \Sigma} d(\phi, \mu_0, \mu_1, \Sigma) \xrightarrow{\text{loss}(1-)} \phi = \frac{\sum_{i=1}^m y^i / m}{\sum_{i=1}^m \mathbb{1}\{y^i = 1\}} / m$$

$$\mu_0 = \frac{\sum_{i=1}^m \mathbb{1}\{y^i = 0\} x^i}{\sum_{i=1}^m \mathbb{1}\{y^i = 0\}}$$



indicator function:  $\mathbb{1}\{\text{true}\} = 1$ ,  $\mathbb{1}\{\text{false}\} = 0$ .

$$\mu_1 = \frac{\sum_{i=1}^m \mathbb{1}\{y^i = 1\} x^i}{\sum_{i=1}^m \mathbb{1}\{y^i = 1\}}$$

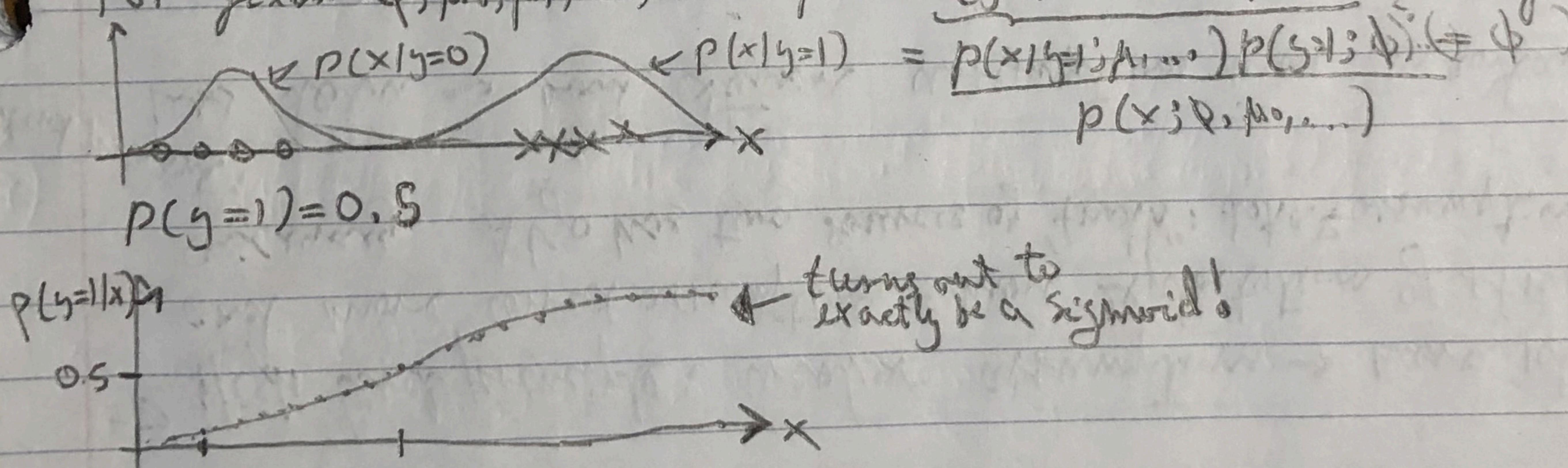
$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^i - \mu_0 y^i)(x^i - \mu_1 y^i)^T$$

Prediction:

$$\arg\max_y P(y|x) = \arg\max_y \frac{p(x|y) p(y)}{p(x)} \\ = \arg\max_y P(x|y) P(y)$$

Compare GDA to logistic regression

For fixed  $\phi, \mu_0, \mu_1, \Sigma$ , let's plot  $p(y=1|x; \phi, \mu_0, \mu_1, \Sigma)$  as a function of  $x$ .



(generative) GDA assumes

$$\begin{aligned} x|y=0 &\sim N(\mu_0, \Sigma) \\ x|y=1 &\sim N(\mu_1, \Sigma) \\ y &\sim \text{Ber}(0.5) \end{aligned}$$

(discriminative) logistic regression

$$p(y=1|x) = \frac{1}{1 + e^{-\phi^T x}} ("x_0=1") \text{ (logistic)}$$

Thus GDA makes stronger assumptions than logistic.

Note: In general, if you make stronger assumptions that are correct, your model will do better (you exploit structure). However, if they aren't, it'll do worse.

$$\left. \begin{array}{l} x|y=1 \sim \text{Poisson}(\lambda_1) \\ x|y=0 \sim \text{Poisson}(\lambda_0) \\ y \sim \text{Pr}(\phi) \end{array} \right\} \Rightarrow p(y=1|x) \text{ is logistic}$$

(or some other  
expn. family)

Thus if you don't know whether your data is Gaussian or Poisson,  
logistic regn. will work regardless.

$\Rightarrow$  High level principle: weaker assumptions  $\rightarrow$  more robust to <sup>assumption</sup> weaker, and more correct assumptions = more exploitation of structure = better model accuracy. Algo has two sources of "truth": data & assumptions.

<sup>Important principle:</sup> Less data  $\rightarrow$  good, extensive assumptions more important.

### Naive Bayes

Feature vector  $x$ ?

Features  
 ↘  
 a word  
 count  
 ↗  
 buy  
 ↘  
 synergy

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ i \\ \vdots \\ 0 \end{bmatrix}, x \in \mathbb{R}^{10,000}$$

$$x_i = \mathbf{1}\{\text{word } i \text{ appears in email}\}$$

Want to model  $p(x|y), P(y)$

~~2<sup>10,000</sup> - 1~~ possible values for  $x$

Assume  $x_i$ 's are conditionally independent given  $y$ .  $p(x_1, \dots, x_{10,000}|y) = p(x_1|y) p(x_2|y) p(x_3|y) \dots p(x_{10,000}|y)$

$$\begin{aligned} &\stackrel{\text{assume}}{=} p(x_1|y) p(x_2|y) p(x_3|y) \dots p(x_{10,000}|y) \\ &= \prod_{i=1}^n p(x_i|y) \end{aligned}$$

↑  
 Naive Bayes assumption  
 Not always great, but sometimes works.

Parameters:

$$\phi_{j|y=1} = p(x_j=1 | y=1), \phi_{j|y=0} = p(x_j=1 | y=0), \phi_y = p(y=1)$$

Joint likelihood:  $L(\phi_s, \phi_{j|y}) = \prod_{i=1}^m p(x_i^i, y_i^i; \phi_y, \phi_{j|y})$ .

MLE:

$$\hat{\phi}_y = \frac{\sum_{i=1}^m \mathbf{1}\{y_i^i = 1\}}{m}$$

$$\hat{\phi}_{j|y=1} = \frac{\sum_{i=1}^m \mathbf{1}\{x_j^i = 1, y_i^i = 1\}}{\sum_{i=1}^m \mathbf{1}\{y_i^i = 1\}}$$

It's not too shabby in practise,.. And extremely efficient.

## Lecture 6

Outline:

Naive Bayes

- Laplace smoothing
- Event Models

Comments on applying ML

- SVM intro

Recap:

$$x = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{bmatrix} \quad \text{a word} \Leftrightarrow x_j = 1 \{ \text{word } j \text{ appears in email} \}$$

Generative model:  $[P(x|y)] [P(y)]$

$$P(x|y) = \prod_{j=1}^n P(x_j|y)$$

$$\text{Parameters: } P(y=1) = \phi_y, P(x_j=1|y=0) = \phi_{j|y=0}, \\ P(x_j=1|y=1) = \phi_{j|y=1}$$

We found:

$$\phi_y = \sum_{i=1}^m I\{y_i=1\} / m$$

$$\phi_{j|y=0} = \sum_{i=1}^m I\{x_i=1, y_i=0\} / \sum_{i=1}^m I\{y_i=0\}$$

At prediction time:

$$P(y=1|x) = \frac{p(x|y=1)P(y=1)}{p(x|y=1)P(y=1) + p(x|y=0)P(y=0)}$$

(majority rule)

$$\text{Example: MBS} \Leftrightarrow j=6017 \Rightarrow P(x_{6017}=1|y=1) = \frac{0}{\#\{y=1\}} = \phi_{6017|y=1}$$

$$P(x_{6017}=1|y=0) = \frac{0}{\#\{y=0\}} = \phi_{6017|y=0} \quad [\text{Haven't seen this word yet.}]$$

$\Rightarrow$  When you suddenly start seeing MBS in your emails,  $P(x|y=1) = \frac{10^{1000}}{1} P(x|y)=0$   
 and so will  $p(x|y=1)$  &  $p(x|y=0)$ , in denominator of  $P(y=1|x) \Rightarrow 0/0$  problem

$\Rightarrow$  Addresses this problem with "Laplace Smoothing". But before that: an example.

Won?

9/12 Valeeforest  
10/10 Oregon State  
10/17 Arizona  
11/21 Caltech  
12/31 Oklahoma

0

0

0

0

?

what's the estimate for this?  
they lost in the end lol!

$$P(x=1) = \frac{\#1's}{\#1's + \#0's} = \frac{0}{0+4} = 0. \quad \text{not a good idea to estimate this.}$$

Laplace smoothing: just add 1 to wins & losses, giving

$$P(x=1) = \frac{1}{1+5} = \frac{1}{6}. \quad \text{Note: This trick is actually optimal in some conditions.}$$

More generally:  $x \in \{1, \dots, k\}$  Estimate  $P(x=j) = \frac{\sum_{i=1}^m \mathbb{1}\{x_i=j\} + 1}{m+k}$

$$P(x=j|y=0) = \left( \sum_{i=1}^m \mathbb{1}\{x_i=1, y_i=0\} + 1 \right) / \left( \sum_{i=1}^m \mathbb{1}\{y_i=0\} + 2 \right) \quad \text{removes 0/0 problem}$$

$x_i \in \{1, \dots, k\}$

rule of thumb is 10 buckets tho, not 4.

size	< 400 size	400-800	800-1200	> 1,200
x	1	2	3	4
1				
2				
3				
4				

$$P(x|y) = \prod_{i=1}^h p(x_j|y)$$

multinomial

"Multivariate Bernoulli event model"

so far:

$x_2$	0	1	2	3	4
0	aardvark				
1		drugs			
2			drugs over		
3				over	
4					over

Multinomial event model

New rep:

$$x \in \begin{bmatrix} 1600 \\ 800 \\ 1600 \\ 8200 \end{bmatrix} \in \mathbb{R}^n \quad x_j = \{1, \dots, 10'000\}$$

$n_i = \text{length of email } i$

$$P(x, y) = P(x|y)P(y) \stackrel{\text{assume } n}{=} \prod_{j=1}^n p(x_j|y)$$

how a "multinomial" instead of "Bernoulli"  
probability  
change of word  $i$  being  $k$  is  
 $y=0$ .

Parameters :  $\phi_{y=1} = P(y=1)$ ,  $\phi_{k|y=0} = \overbrace{P(x_j=k|y=0)}^{\text{assume independent of } j \text{ (location)}}$

$$\phi_{k|y=1} = p(x_i=k|y=1)$$

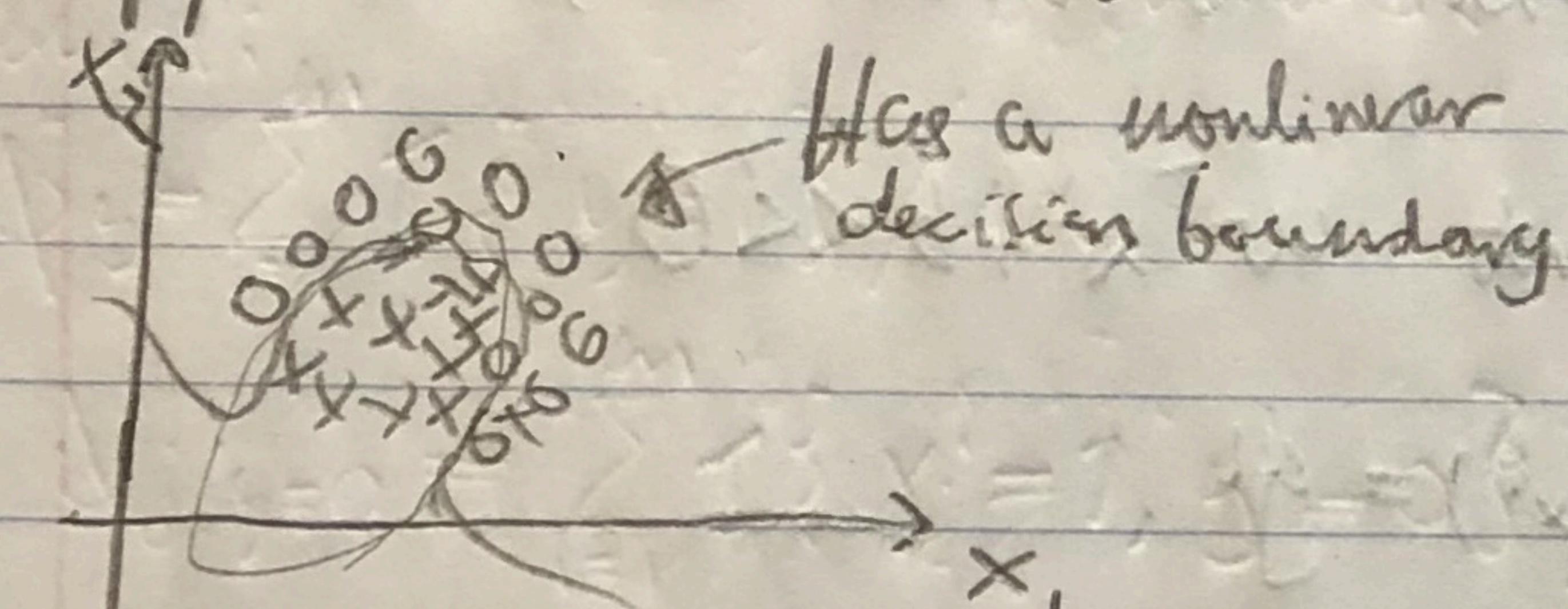
$$\text{MLE : } \phi_{k|y=0} = \frac{\left( \sum_{i=1}^m \mathbb{1}\{y^i=0\} \sum_{j=1}^n \mathbb{1}\{x_j^i=k\} \right) / m}{\left( \sum_{i=1}^m \mathbb{1}\{y^i=0\} \cdot n_i + 10,000 \right)}$$

number of possible outcomes

Naive Bayes is Logistic regression in practise :

Log Reg typically performs better, but Naive Bayes much simpler to implement and far more computationally efficient (no need for eg grad. Descent). Naive Bayes is a nice "quick & dirty" starting point when investigating data.

## Support Vector Machines



SVM maps to higher dimension features:  $\phi(x) = \begin{bmatrix} x_1 \\ x_2 \\ x_1x_1 \\ x_1x_2 \\ x_2x_2 \end{bmatrix}$

- Optimal margin classifier (separable case)
- Kernels:  $x \in \mathbb{R}^2 \rightarrow \phi(x) \in \mathbb{R}^{1000}$  or even  $\mathbb{R}^\infty$
- Inseparable case

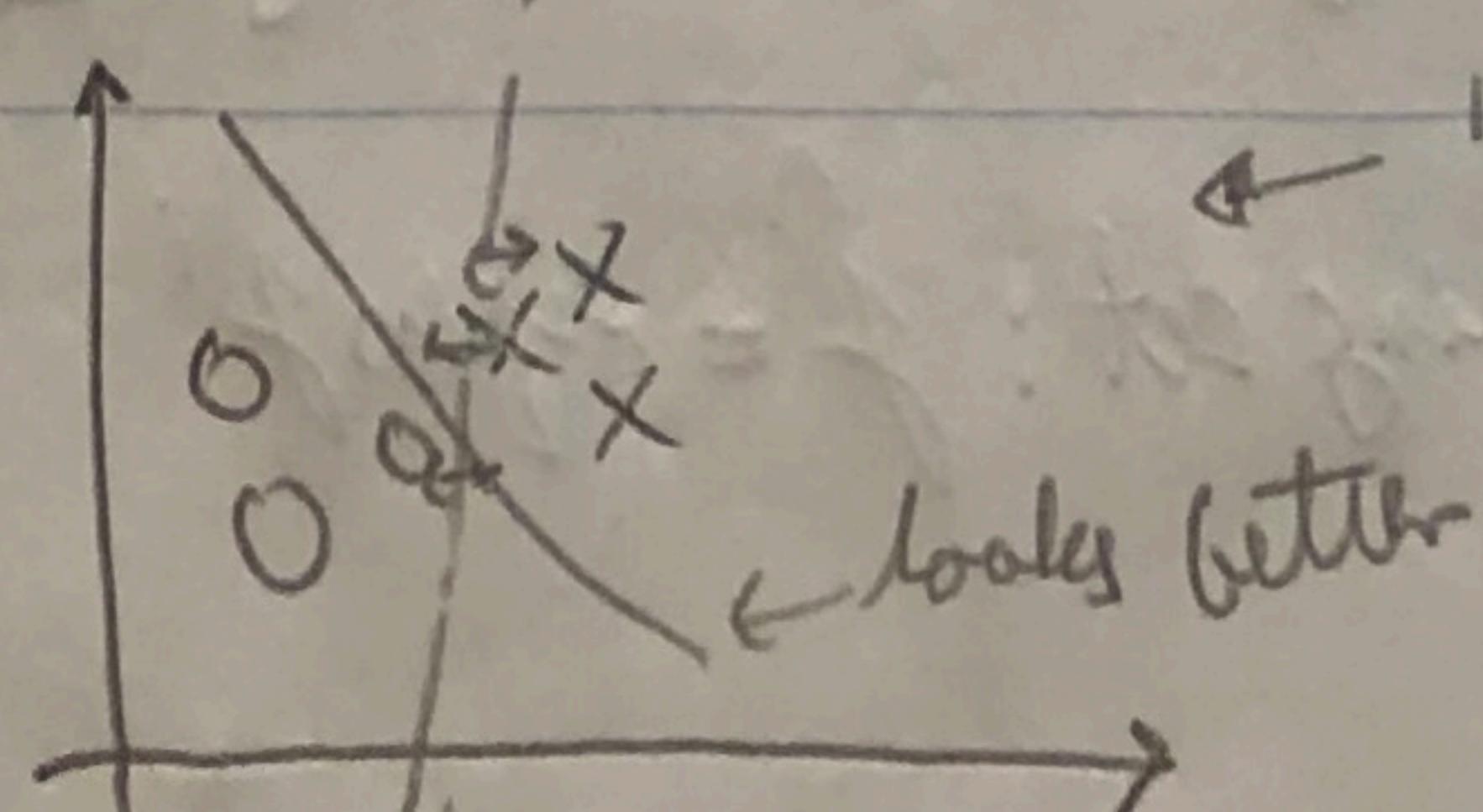
## Functional Margin

$$h_\theta(x) = g(\theta^T x)$$

Predict "1" if  $\theta^T x > 0$  ( $g(\theta^T x) \geq 0.5$ )  
≠ 0 otherwise

If  $y^i=1$ , hope that  $\theta^T x^i > 0$  & if  $y^i=0$ , hope  $\theta^T x^i < 0$

Geometric Margin:



"Opt. Marg. Class." want to find the best line, i.e. geometric margin optimised

Notation:

Labels  $y \in \{-1, +1\}$ . Have  $h$  output values in  $\{-1, +1\}$

$$g(z) = \begin{cases} 1 & \text{if } z \geq 0, \\ -1 & \text{otherwise} \end{cases}$$

Previously:  $\mathbb{R}^{n+1}, x_0 = 1$

$$h_{w,b}(x) = g(w^T x)$$

Now:

$$h_{w,b}(x) = g(\underbrace{w^T x + b}_{\sum_{i=1}^n w_i x_i + b}) \quad \text{Drop the } x_0 = 1 \text{ constraint.} \quad \text{Eg: } \left[ \begin{array}{c} b_0 \\ b_1 \\ b_2 \\ b_3 \end{array} \right] \overset{\exists b}{\brace} w$$

Functional margin of hyperplane defined by  $(w, b)$  wrt  $(x^i, y^i)$ :

$$\hat{y}^i = y^i (w^T x^i + b). \quad \text{If } y^i = 1, \text{ want } w^T x^i + b \gg 0 \quad \&$$

$$\text{If } y^i = -1, \text{ want } w^T x^i + b \ll 0.$$

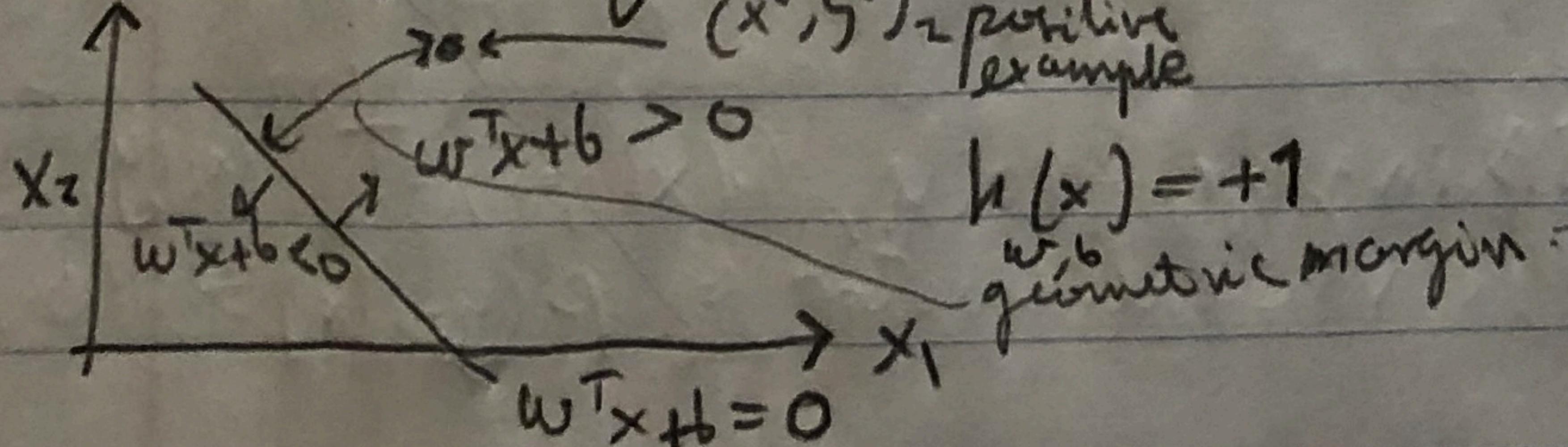
$$\Rightarrow \text{Want } \hat{y}^i \gg 0. \quad \text{If } \hat{y}^i > 0, \text{ then } h(x^i) = y^i$$

Functional margin wrt training set:  $\hat{f} = \min_i \hat{y}^i, i=1, \dots, m$

i.e. it's defined by the "worst" training example

$$\|w\|=1, (w, b) \rightarrow \left( \frac{w}{\|w\|}, \frac{b}{\|w\|} \right) \quad \text{Nice normalisation convention}$$

Geometric Margin



Geometric margin of hyperplane  $(w, b)$  wrt  $(x^i, y^i)$ :  $\hat{y}^i = \frac{w^T x^i + b}{\|w\|}$

just Euclidean distance

$$\text{Thus out: } f' = \frac{\hat{f}'}{\|w\|} \quad \left[ \text{note: } \hat{f}' = \text{functional margin} \right] \quad \left[ \text{f' = geometric margin} \right]$$

Geometric margin wrt training set:  $\hat{f} = \min_i \hat{y}^i$

Optimal Margin classifier;  
Choose  $w, b$  to maximize  $\gamma$ .

$$\hat{\gamma} = \max_{\gamma, w, b} \gamma$$

$$\text{s.t. } y_i(w^T x^i + b) / \|w\| \geq \gamma, i=1, \dots, m$$

nonconvex

This is convex

$$\min \|w\|^2$$

$$\text{w.r.t. } b$$

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

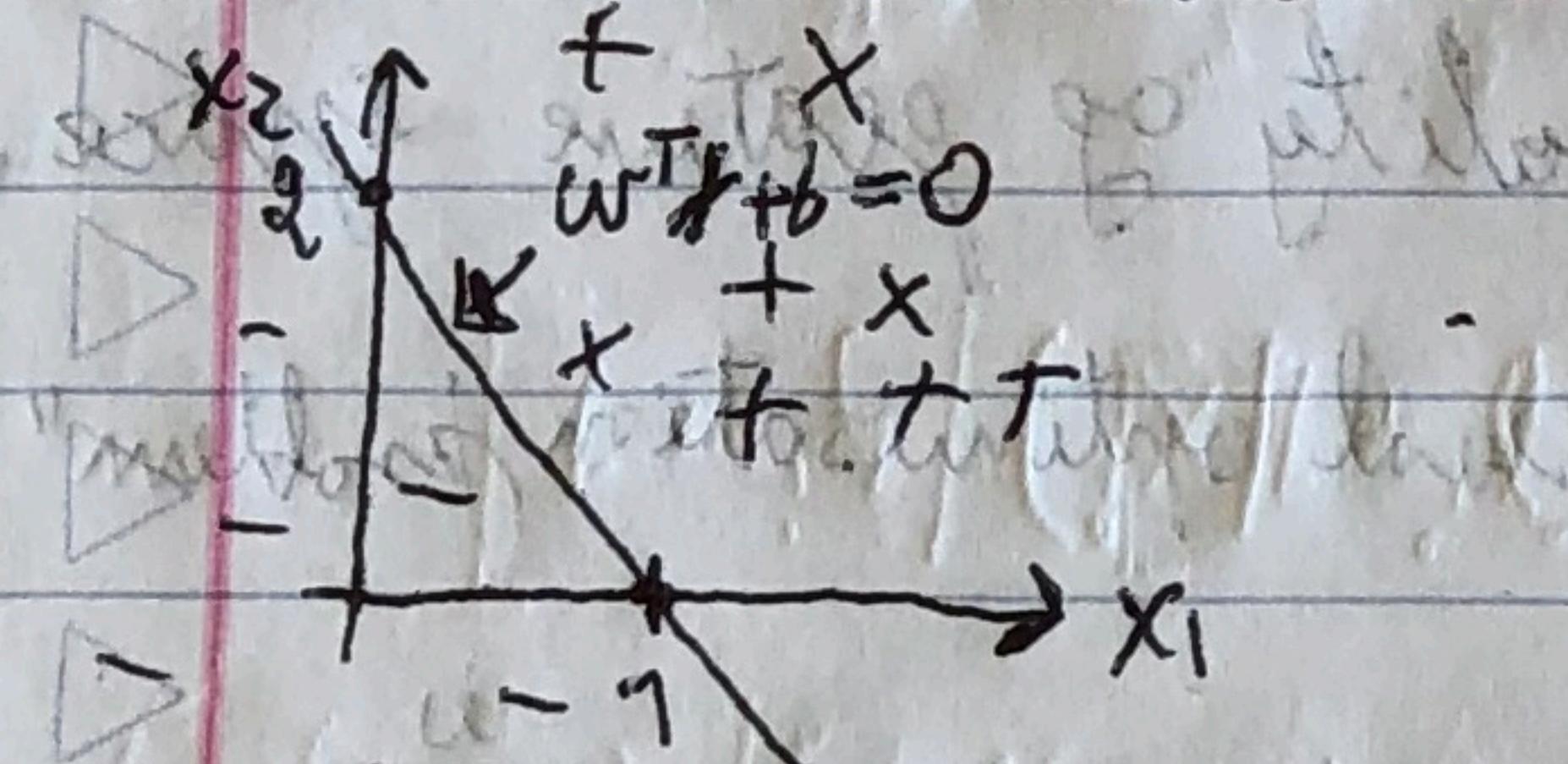
example

Lecture 7: SVMs, Optimisation problem; Representer theorem; kernels.

$$f^i = \frac{y^i(w^T x^i + b)}{\|w\|} \quad (\text{geometric margin}) = ; f = \min_{i=1, \dots, n} f^i \quad ; y = g(w^T x + b) \quad \text{classifier}$$

$$\max_{w, b} f \text{ st } \frac{y^i(w^T x^i + b)}{\|w\|} \geq \gamma \quad \text{equiv. to "every example has GM} \geq \gamma \text{"}$$

Note: Can scale  $w, b$  without affecting decision boundary. Eg, take  $y = g([z]^\top x - z) = g([x_1 \ x_2]^\top x - 2)$



$\Rightarrow$  Optimisation objective becomes  $\max_{w, b} \frac{1}{\|w\|} \text{ st } y^i(w^T x^i + b) \geq \gamma$

$$\Leftrightarrow \min_{w, b} \frac{1}{2} \|w\|^2 \text{ st } y^i(w^T x^i + b) \geq 1, i=1, \dots, n$$

Say  $x^i \in \mathbb{R}^{100}$ . Suppose  $w = \sum_{i=1}^n \alpha_i x^i$ . The "Representer theorem" proves this to be an optimal value, but proof is quite involved. Here's an intuitive explanation:

Intuition: Logistic regression:  $\theta_0 = 0$ , then  $\theta := \theta - \alpha(h_\theta(x^i) - y^i)x^i$ . Thus  $\theta$  is a linear combination of the training examples.  $\Rightarrow \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \leftarrow \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_n \end{bmatrix}$

Also true for batch gradient descent:  $\theta := \theta - \alpha \sum_{i=1}^n (h_\theta(x^i) - y^i)x^i$

\*2: Eg:  $g([z]^\top x - z) \Rightarrow w$  is always normal to decision boundary. Then  $w$  lies in the span of the training examples:

\*3:  $w$  is correctly placed in the span of the training example, i.e. it's safely represented in the features  $x_1, x_2$ , where the training example is.

So let's assume  $w = \sum_{i=1}^n \alpha_i y^i x^i$ .  $\min_{w, b} \frac{1}{2} \|w\|^2$  st  $y^i (w^T x^i + b) \geq 1 \forall i$ .

$$\Leftrightarrow \min_{w, b} \frac{1}{2} \left( \sum_{i=1}^n \alpha_i y^i x^i \right)^T \left( \sum_{j=1}^n \alpha_j y^j x^j \right) = \min_{w, b} \frac{1}{2} \sum_{i,j} \sum_{i,j} \alpha_i \alpha_j y^i y^j x^i x^j = \frac{1}{2} w^T w$$

$\langle x, z \rangle = x^T z$  is the inner product

while the limitation becomes:  $y^i (\sum_j \alpha_j \langle x^j, x^i \rangle + b) \geq 1 \Rightarrow$  feature vectors only

ever appear in inner product, and so we only need to be able to evaluate it very efficiently, regardless of the dimensionality of feature space.

Anyways, can simplify that optimisation problem to:

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \sum_{i,j} y^i y^j \alpha_i \alpha_j \langle x^i, x^j \rangle \text{ st } \alpha_i \geq 0, \sum_{i=1}^n y^i \alpha_i = 0$$

"Dual optimisation problem"

To predict: 1) Solve for  $\alpha^*, b$ ; 2) Predict by:  $h_{w,b}(x) = g(w^T x + b) = g(\underbrace{(\sum_i \alpha^* y^i x^i)^T}_{=w} x + b)$   
 $= g(\sum_i \alpha^* y^i \langle x, x^i \rangle + b)$ .

Kernel trick:

1) Write algorithms in terms of  $\langle x^i, x^j \rangle$  ( $\langle x, z \rangle$ )

2) Let there be some mapping from  $x \mapsto \phi(x)$

3) Find way to compute  $K(x, z) = \phi(x)^T \phi(z)$

4) Replace  $\langle x, z \rangle$  in algorithm with  $K(x, z)$ .

Let  $x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \in \mathbb{R}^n$ ,  $\phi(x) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_1 x_3 \\ x_2 x_1 \\ x_2 x_2 \\ x_2 x_3 \\ x_3 x_1 \\ x_3 x_2 \\ x_3 x_3 \end{bmatrix} \in \mathbb{R}^{n^2}$ ,  $\phi(z) = \begin{bmatrix} z_1 z_1 \\ z_1 z_2 \\ z_1 z_3 \\ z_2 z_1 \\ z_2 z_2 \\ z_2 z_3 \\ z_3 z_1 \\ z_3 z_2 \\ z_3 z_3 \end{bmatrix}$

$n^2$  elements  $\Rightarrow$  need  $O(n^2)$  time to compute  $\phi(x)$  or  $\phi(x)^T \phi(z)$  explicitly.

$k(x, z) = \phi(x)^T \phi(z) = (x^T z)^2$  is an  $O(n)$  time computation.

$$= (\sum_{i=1}^n x_i z_i) (\sum_{j=1}^n x_j z_j) = \sum_{i,j} x_i z_i x_j z_j = \sum_{i,j} (x_i x_j) (z_i z_j) \Rightarrow \text{Go from } O(n^2) \text{ to just } O(n)!$$

and ditto for  $\phi(z)$

Another ex:  $k(x, z) = (x^T z + c)^d$ ,  $c \in \mathbb{R} \Rightarrow \phi(x) \rightarrow$

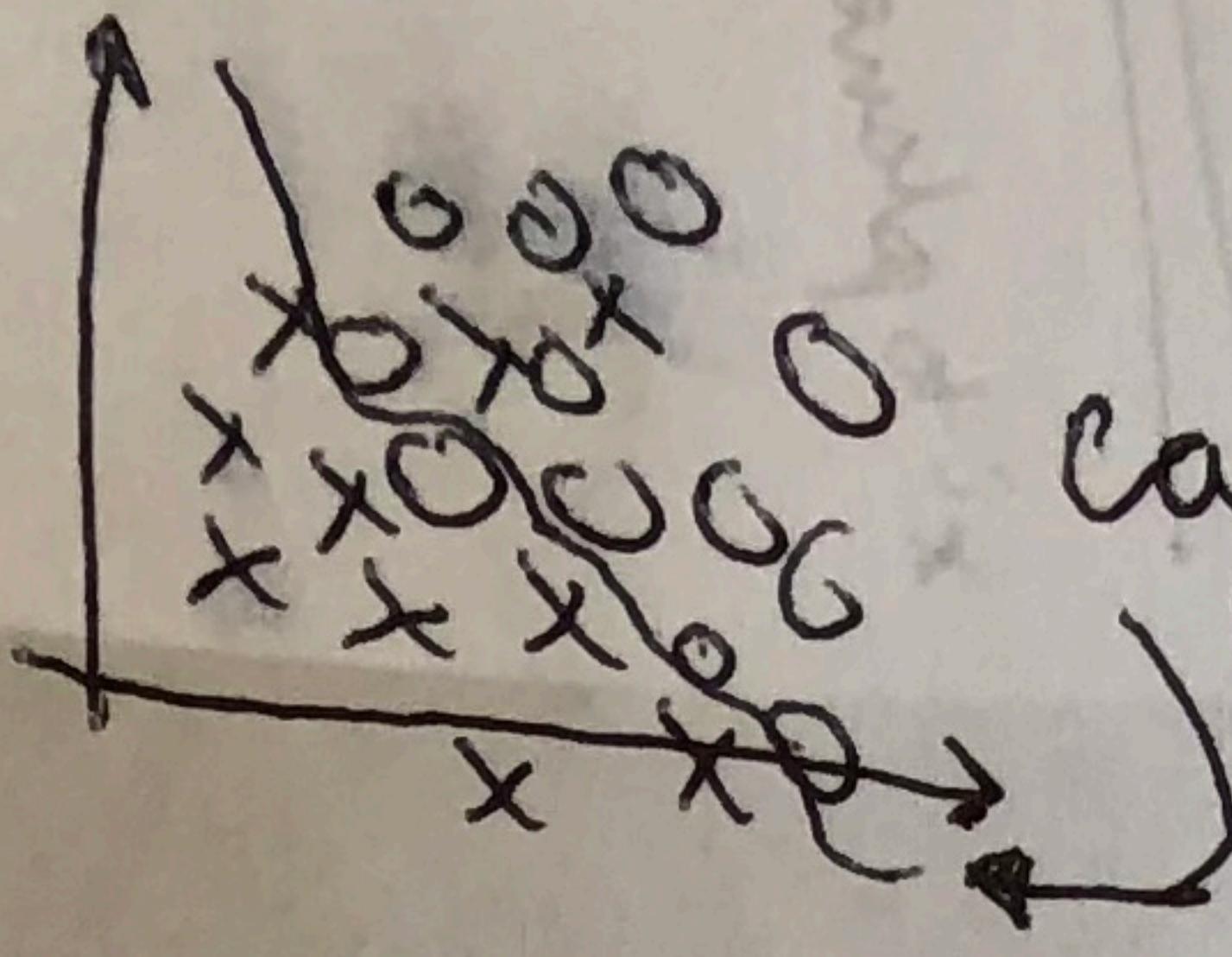
Or:  $k(x, z) = (x^T z + c)^d \Rightarrow \phi(x)$  has all

$\binom{n+d}{d}$  features of monomials up to order  $d$ .

$$\frac{x_1 x_2 \dots x_n x_d}{x_1 x_2 \dots x_n} (n+d)^d \text{ number of features}$$

30

Optimal margin classifier + kernel trick  
= SVM.



can learn very nonlinear decision boundaries.

How to make kernels. Heuristic:

If  $x, z$  are "similar", then  $k(x, z) = \phi^T(x)\phi(z)$  is "large". Else if  $x, z$  are "dissimilar",  $k(x, z)$  should be small. i.e.,  $\phi(x)$  is  $\phi(z)$

$$k(x, z) = \exp\left(-\frac{|x-z|^2}{2\sigma^2}\right), \text{ i.e.}$$

Does there exist  $\phi$  s.t  $k(x, z) = \phi^T(x)\phi(z)$ ? This is the goal. Must satisfy:  $k(x, x) \geq 0$ .

Theorem outlining when a given kernel is appropriate:

Let  $\{x^1, \dots, x^d\}$  be  $d$  points, and let

$K \in \mathbb{R}^{d \times d}$  "kernel/Gram matrix"

$k_{ij} = k(x^i, x^j)$ . Given any vector  $z$ ,

$$z^T K z = \sum_{i,j} z_i k_{ij} z_j = \sum_{i,j} z_i \phi(x^i)^T \phi(x^j) z_j$$

$$= \sum_{i,j} z_i \sum_k [\phi(x^i)]_k [\phi(x^j)]_k z_j = \sum_{i,j} z_i [\phi(x^i)]_k z_j$$

$$[\phi(x^i)]_k z_j = \sum_k \left( \sum_i z_i \phi(x^i)_k \right) z_j \geq 0 \Rightarrow K \text{ is psd}$$

Theorem (Mercer):  $K$  is a valid kernel function (i.e.  $\exists \phi$  s.t  $k(x, z) = \phi(x)^T \phi(z)$ ) iff for any  $d$  points  $\{x^1, \dots, x^d\}$ , the corresponding kernel matrix  $K$  is psd. This agrees with the theorem. The above algebra

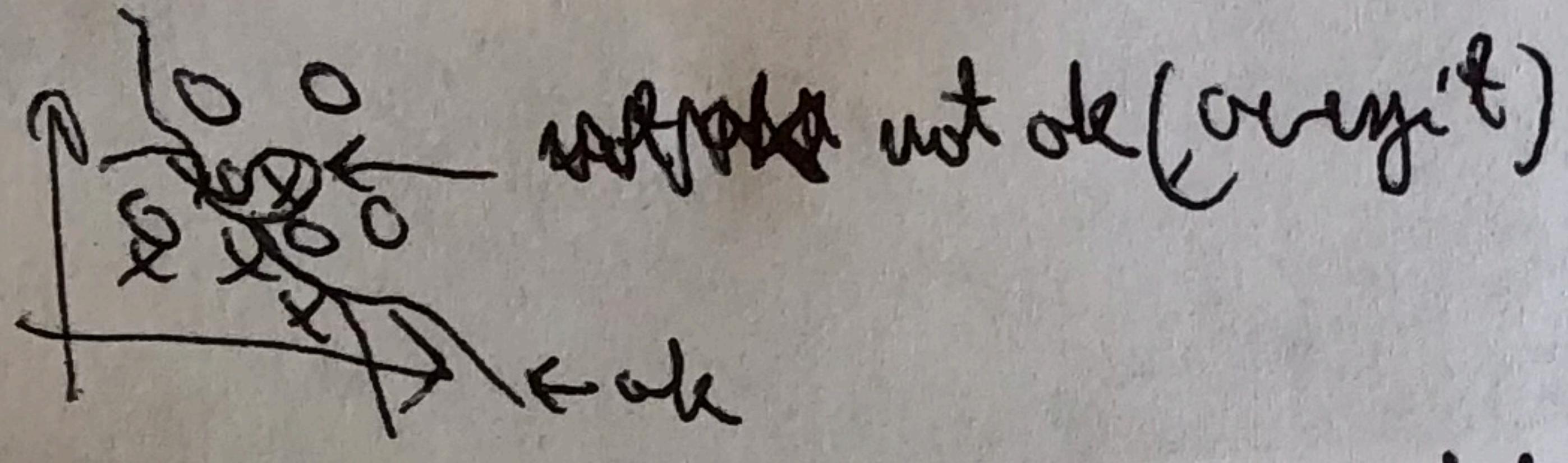
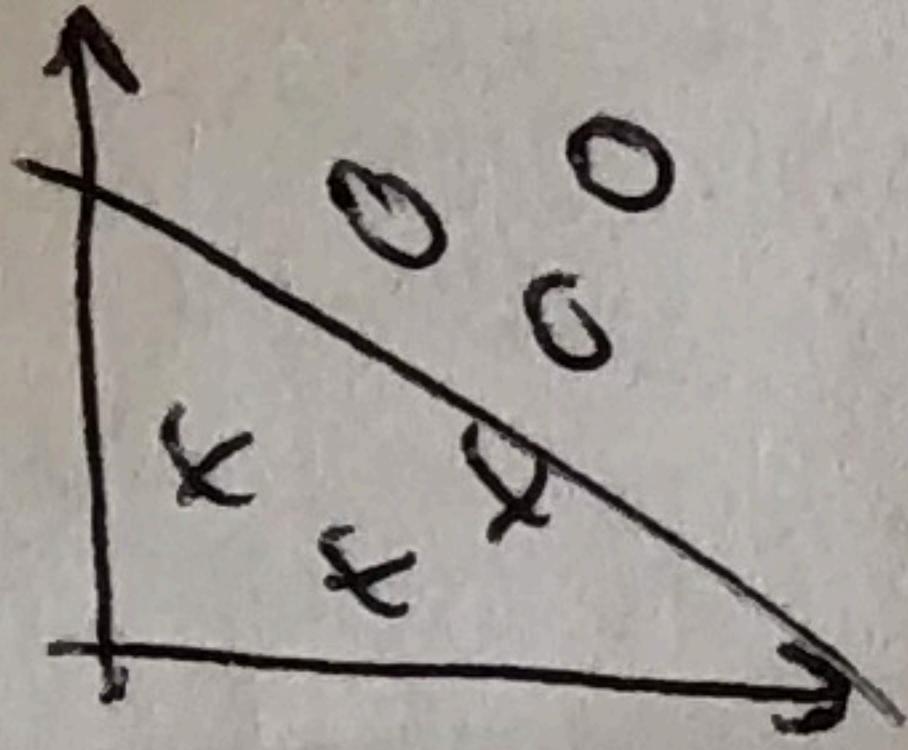
$k(x, z) = \exp\left(-\frac{|x-z|^2}{2\sigma^2}\right)$ , Gaussian kernel, does indeed turn out to be a valid kernel (and is commonly used).

(5)

Most widely used kernel is the linear kernel.  
 $k(x, z) = x^T z$ ;  $\phi(x) = x$ . (Trivial).

After that, the Gaussian kernel is the most widely used, whose  $\phi(x) \in \mathbb{R}^\infty$ . Uses all monomial features:  $x_1, x_1x_2, x_1^2x_2, \dots$  etc.

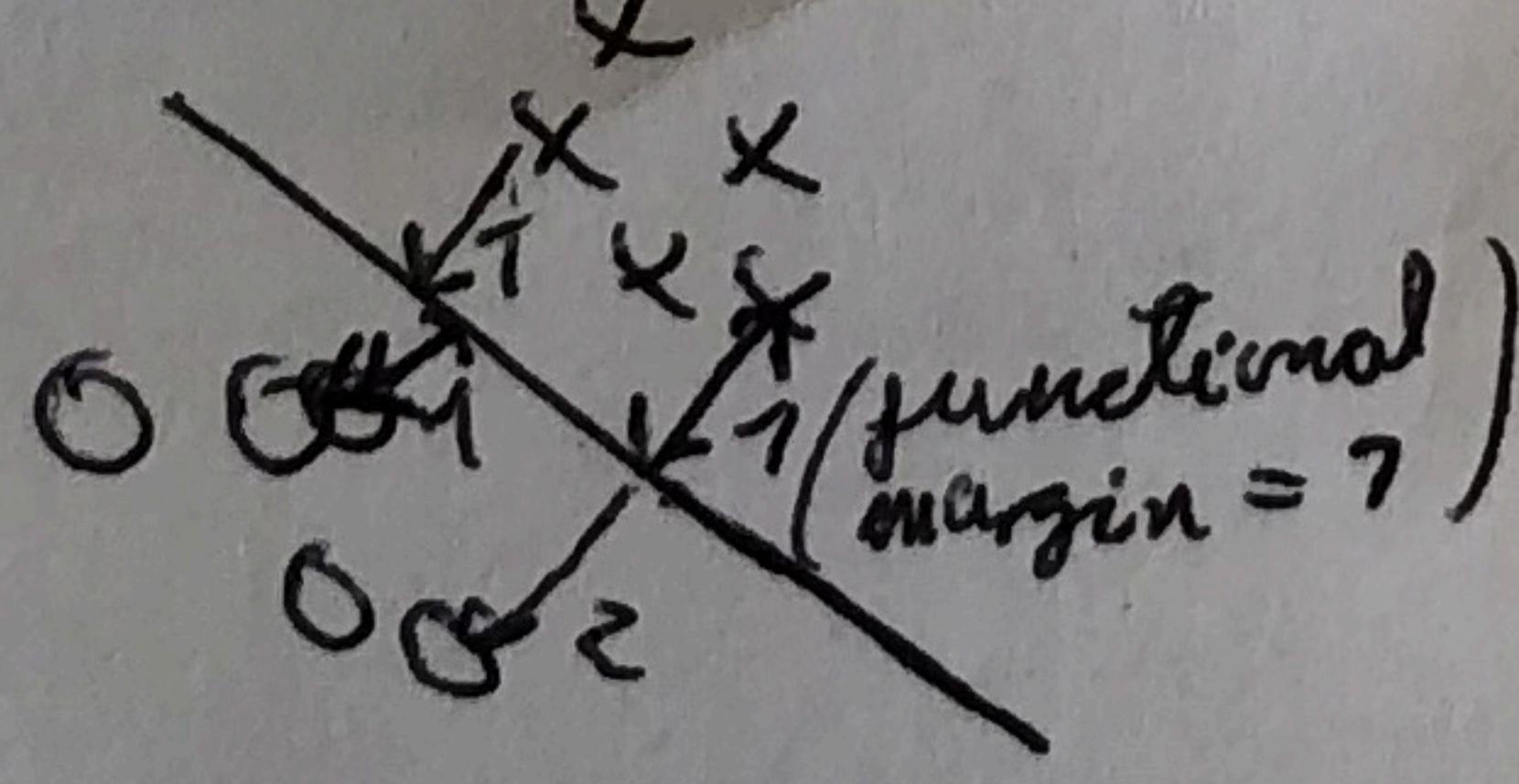
Basically, allows us to do linear regression in infinite dimensional feature spaces.



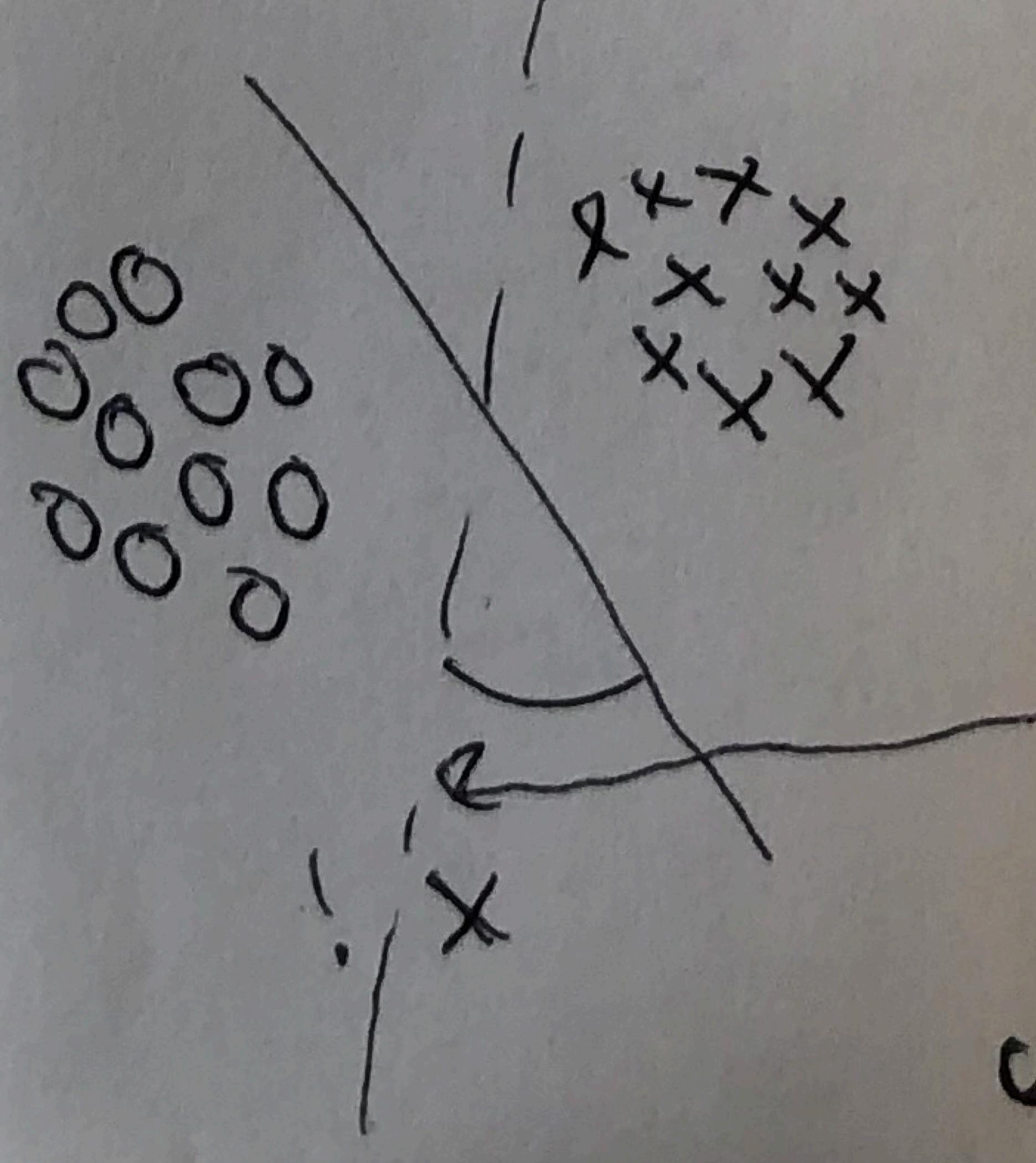
$L_1$  norm soft margin SVM

$$\min \frac{1}{2} \|w\|^2 \text{ st } \underbrace{y_i(w^T x^i + b)}_{\substack{\text{geometric margin} \\ \text{functional}}} \geq 1, \quad i=1, \dots, n$$

$$+ \sum_{i=1}^n \xi_i \quad \text{where } \xi_i \geq 0$$



"Setting  $\xi_i$  to eg 0.5 ~~lets~~ you get away with having a functional margin a bit less than 1".



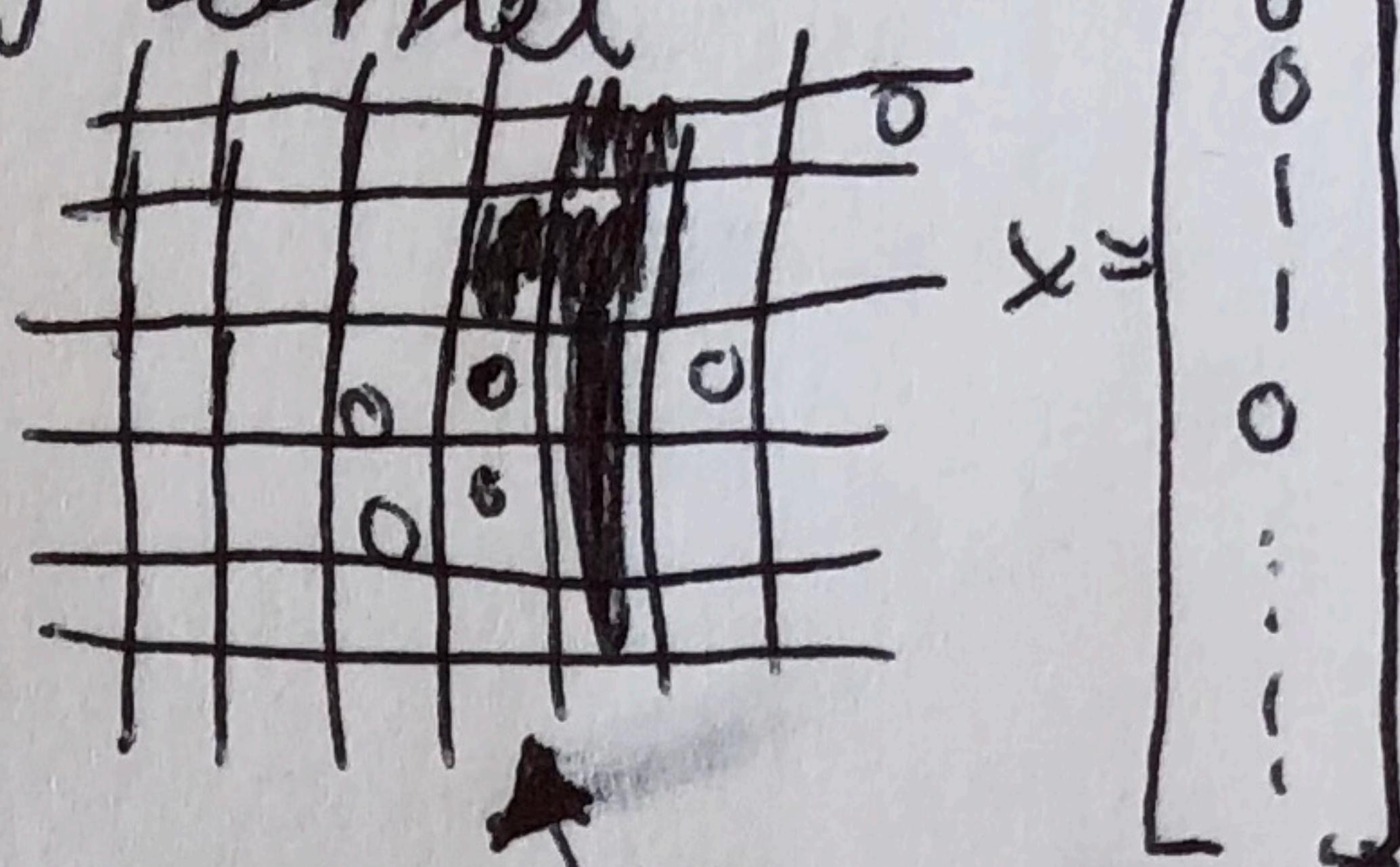
auto margin classifier optimises for the worst case margin, and so will cause a dramatic swing in the decision boundary even just due to a single ~~outlier~~ outlier.  
 Soft margin SVM is less vulnerable to this.

~~$$\max \sum_{i=1}^n \alpha_i - \sum_{i=1}^n \sum_{j=1}^n y_i y_j \alpha_i \alpha_j \langle x^i, x^j \rangle \text{ st } \sum_{i=1}^n y_i \alpha_i = 0,$$~~

$$0 \leq \alpha_i \leq C, i=1, \dots, n$$

$$k(x, z) = (x^T, z)^d \quad \text{poly kernel}$$

$$= \exp\left(-\frac{\|x - z\|^2}{2\sigma^2}\right)$$



$$(7) \quad x = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ -1 \\ 0 \\ \vdots \end{bmatrix}$$

Protein sequence classifier

A...Z

B A J T S T A J B A J T A U ...

$$\phi(x) = ? \quad x \mapsto y$$

AAAA	0
AAAB	0
AAAC	0
AAAD	0
⋮	⋮
AAA Z	2
⋮	⋮
BA JT	2
TS JA	1
ZZZZ	0

$$= \phi(x),$$

$20^4$  length

$$\underline{\phi(x)^T \phi(z) = k(x, z)}$$

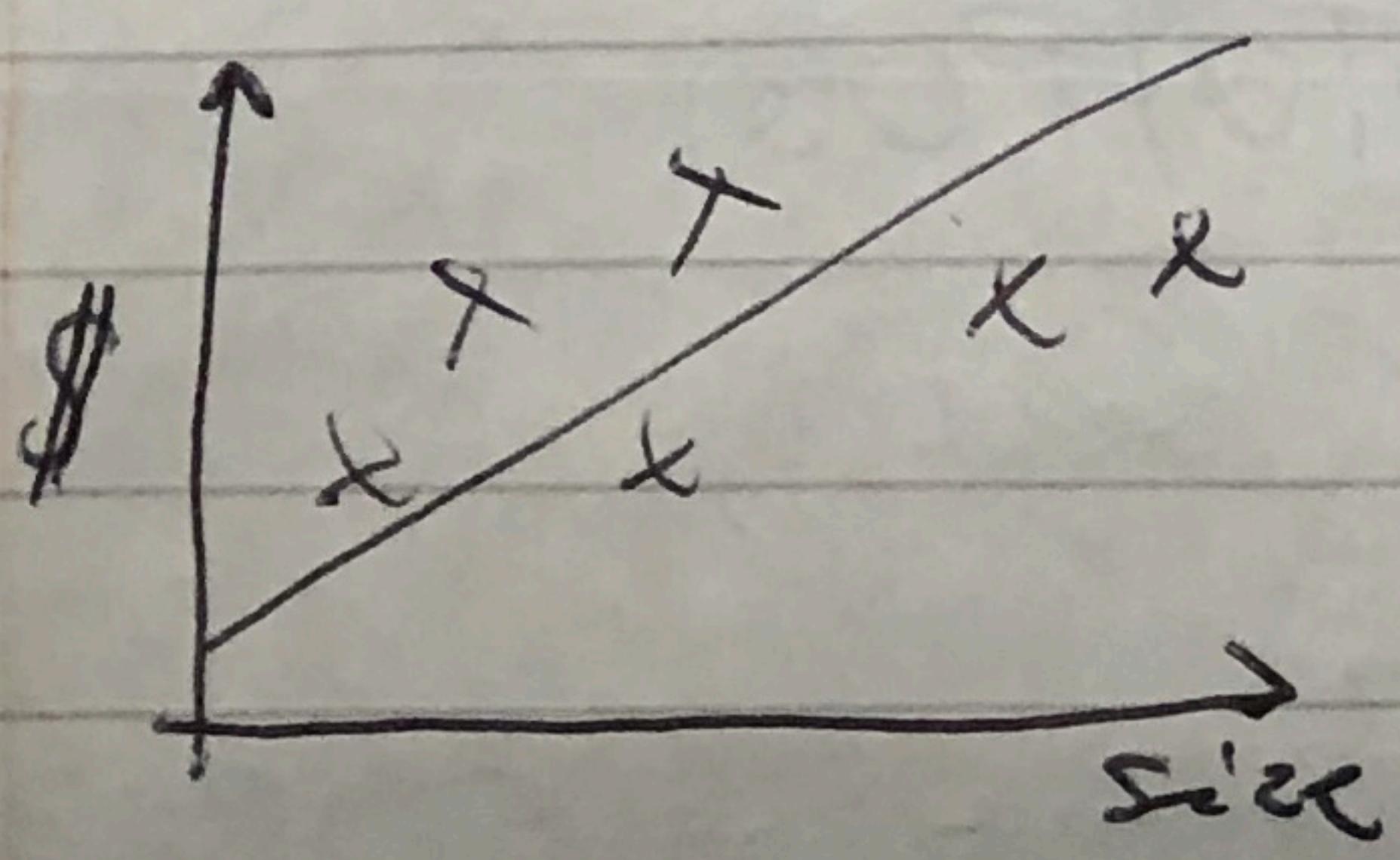
There is a dynamic programming algo to compute this Kernel.

## Lecture 8

### Outline

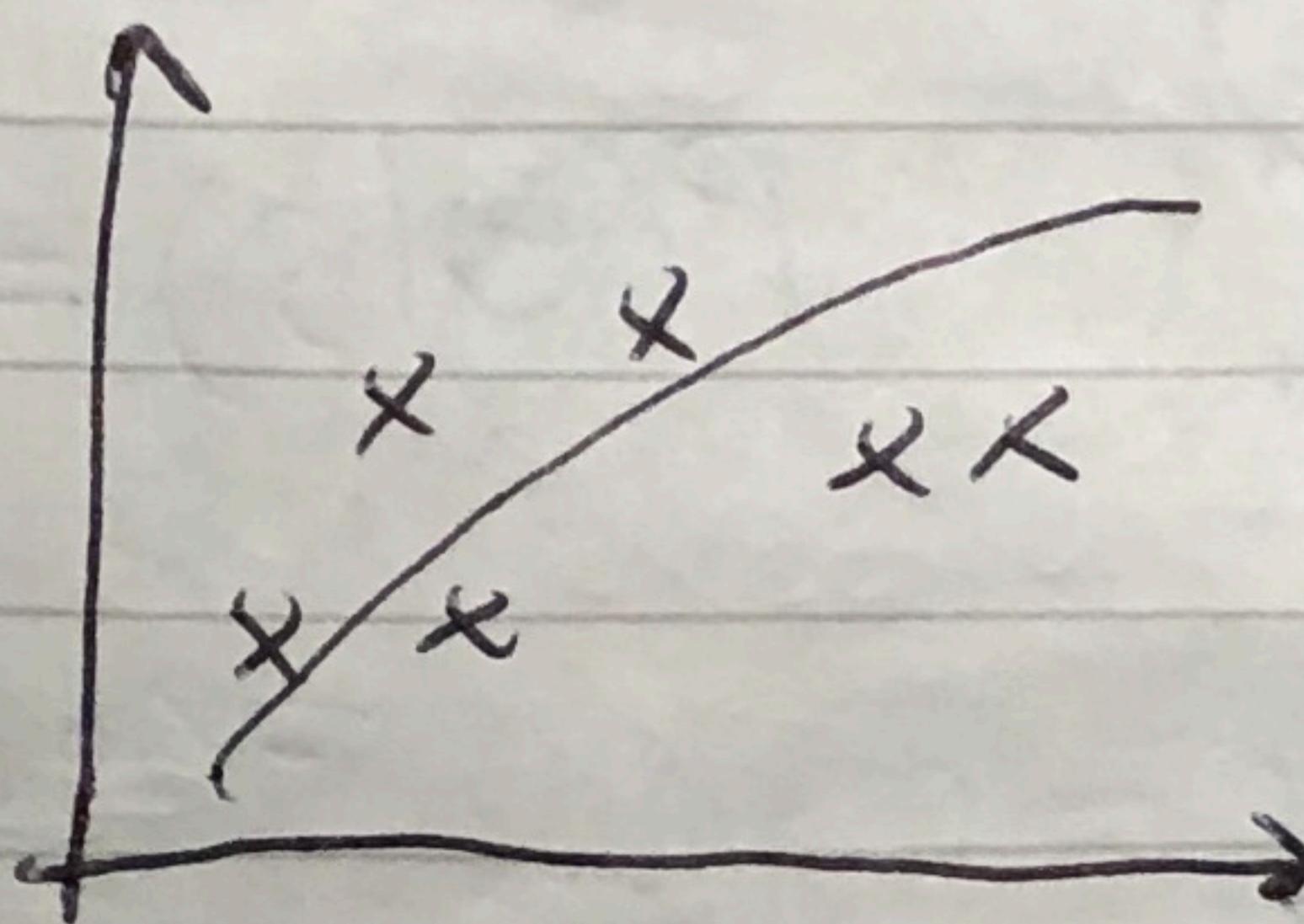
- Bias/variance
- Regularisation
- Train/dev/test splits
- Model Selection & Cross-validation

Bias/variance:



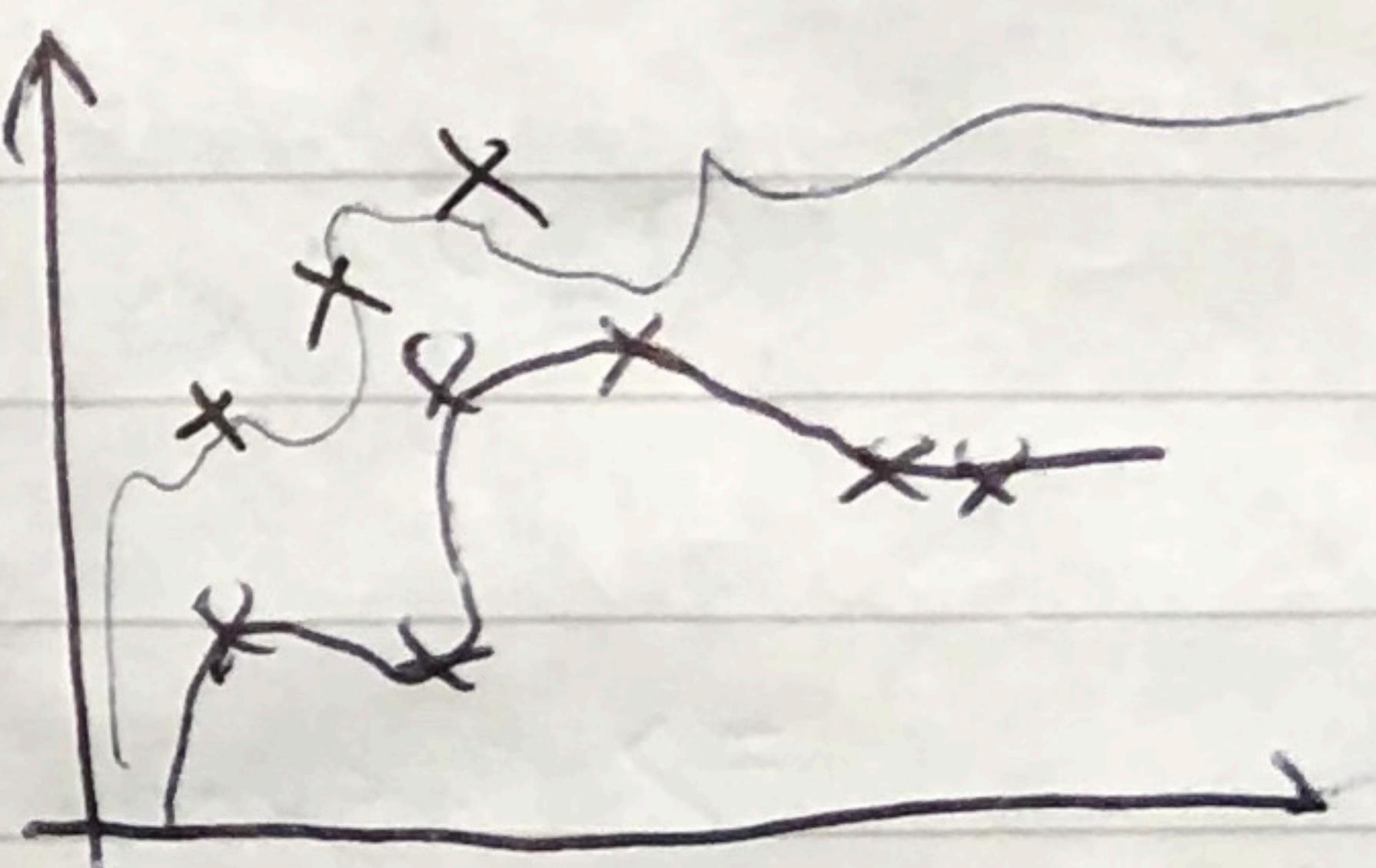
$$\theta_0 + \theta_1 x$$

"Underfit"; high bias



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

"just right"

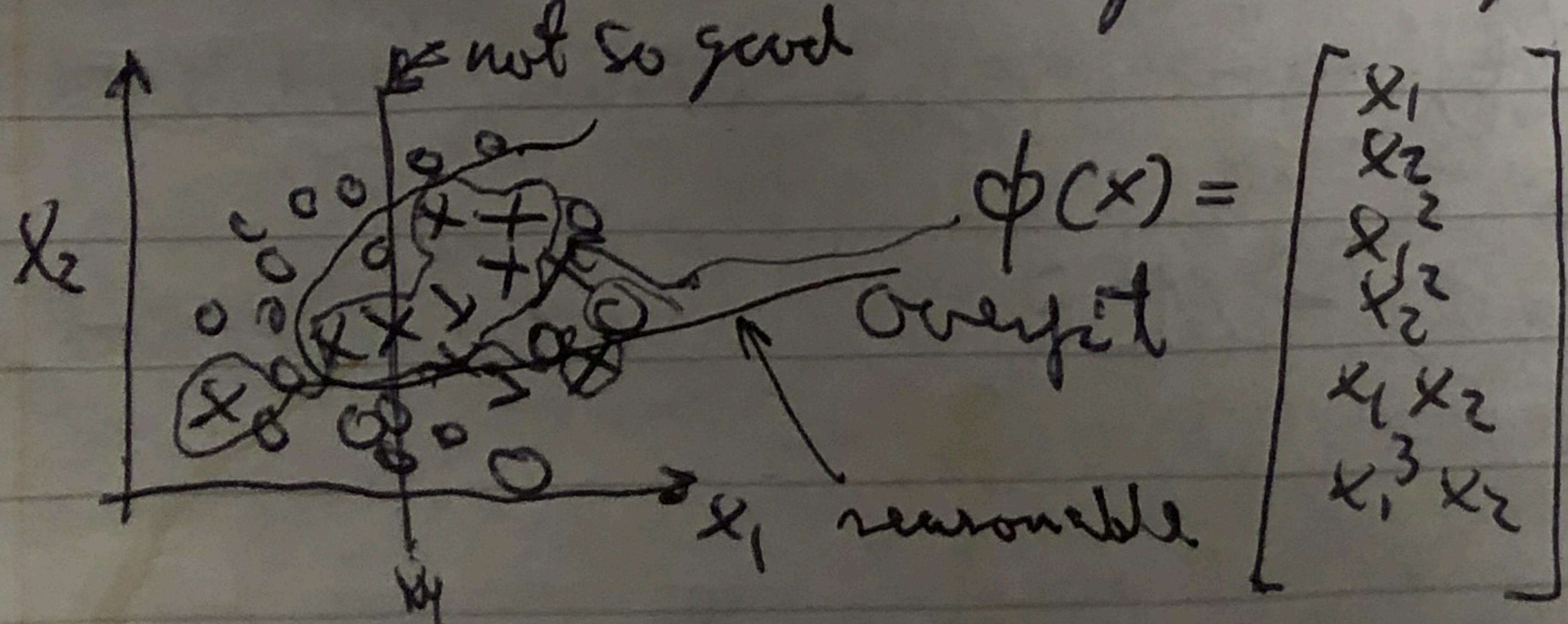


$$\theta_0 + \dots + \theta_5 x^5$$

overfit; high variance  
For new data, this  
algorithm would get  
a completely different  
result!

Bias = "Strong preconception of model". Want to avoid

Also hold true in classification problems



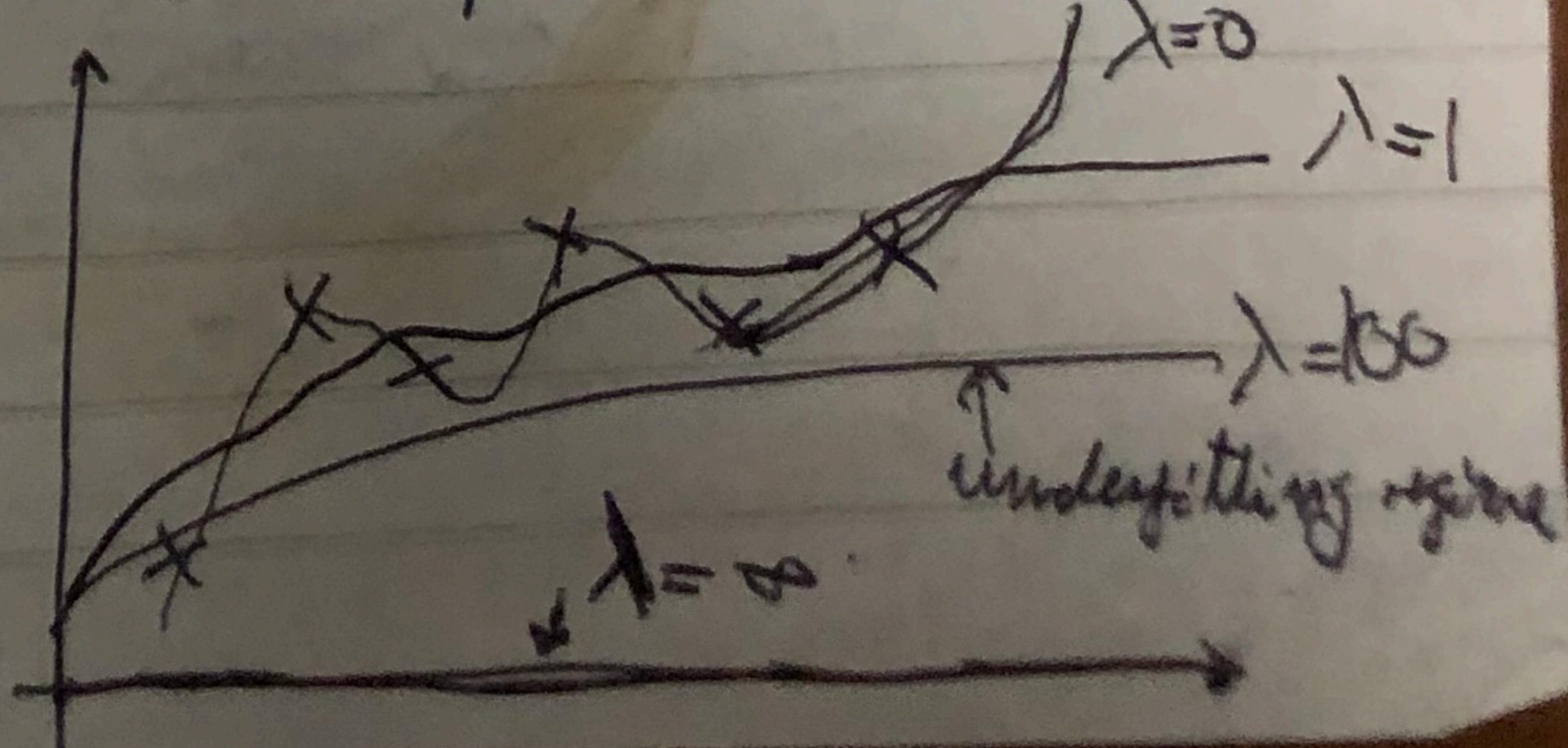
$$\phi(x) = \begin{bmatrix} x_1 \\ x_2 \\ x_1^2 \\ x_2^2 \\ x_1 x_2 \\ x_1^3 x_2 \end{bmatrix}$$

Regularisation: Sounds deceptively simple, but extremely widely used.

$$\min_{\theta} \frac{1}{2} \sum_{i=1}^m |y_i - \theta^T x^i|^2 + \frac{\lambda}{2} |\theta|^2$$

regular. term

removing large  $\theta_0$  term  
& makes overfitting harder



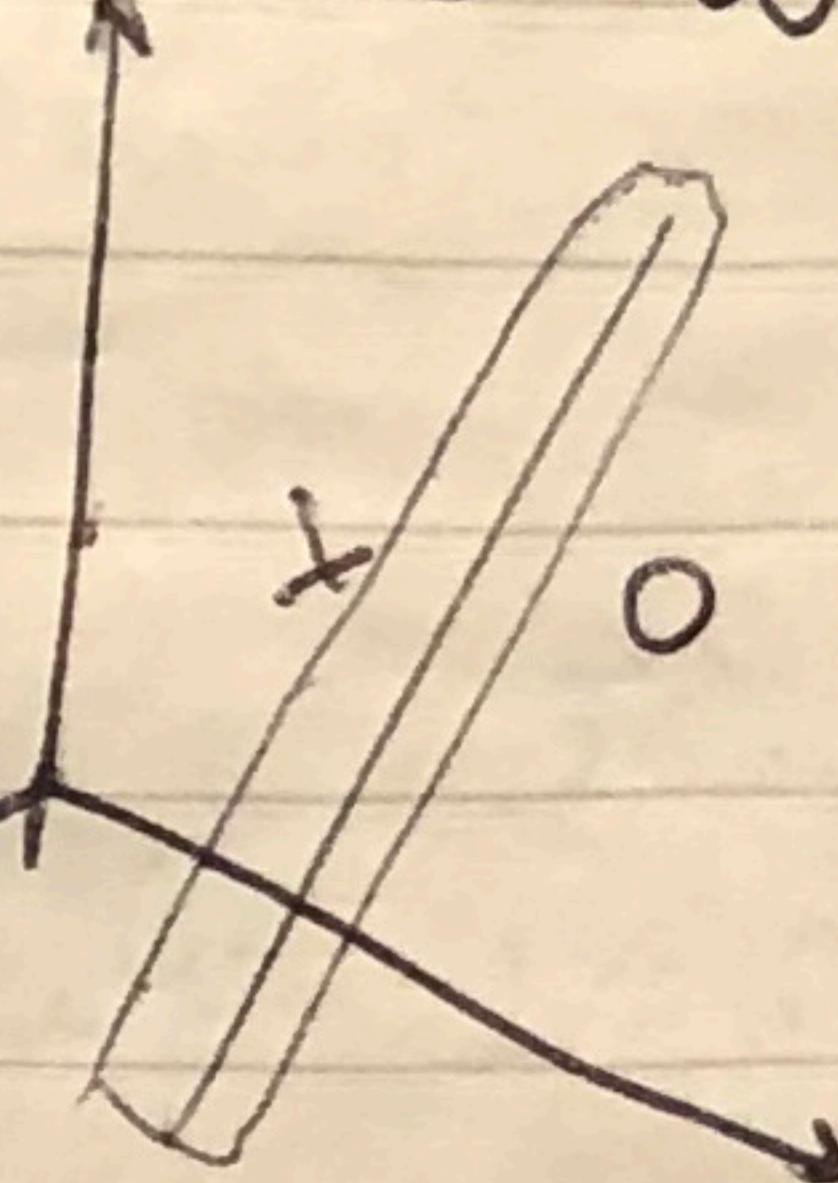
$$\underset{\Theta}{\operatorname{argmax}} \sum_{i=1}^n \log p(y_i | x_i; \Theta) - \frac{\lambda}{2} \|\Theta\|^2$$

Note: SVM don't overfit too much despite working in extremely large feature spaces because  $\|w\|^2 \geq 1/\lambda$ , but formal proof quite complicated.

Text classification:  
 $m = 100$ ;  $n = 10,000$ ,  $X = \begin{bmatrix} 1 & \text{standard} \\ 0 & \\ \vdots & \\ 0 & \\ \vdots & \end{bmatrix}$   
 $\uparrow$  examples     $\uparrow$  features

Need regularisation here to avoid (severe) overfitting.

Note: always ensure features are scaled similarly during preprocessing before applying regularisation.



$$S = \{(x_i, y_i)\}_{i=1}^m \quad P(\Theta|S) = \frac{P(S|\Theta)P(\Theta)}{P(S)}$$

$$\Rightarrow \underset{\Theta}{\operatorname{argmax}} P(\Theta|S) = \underset{\Theta}{\operatorname{argmax}} P(S|\Theta)P(\Theta)$$

$$= \underset{\Theta}{\operatorname{argmax}} \left( \prod_{i=1}^m p(y_i|x_i; \Theta) P(\Theta) \right)$$

assume Gaussian

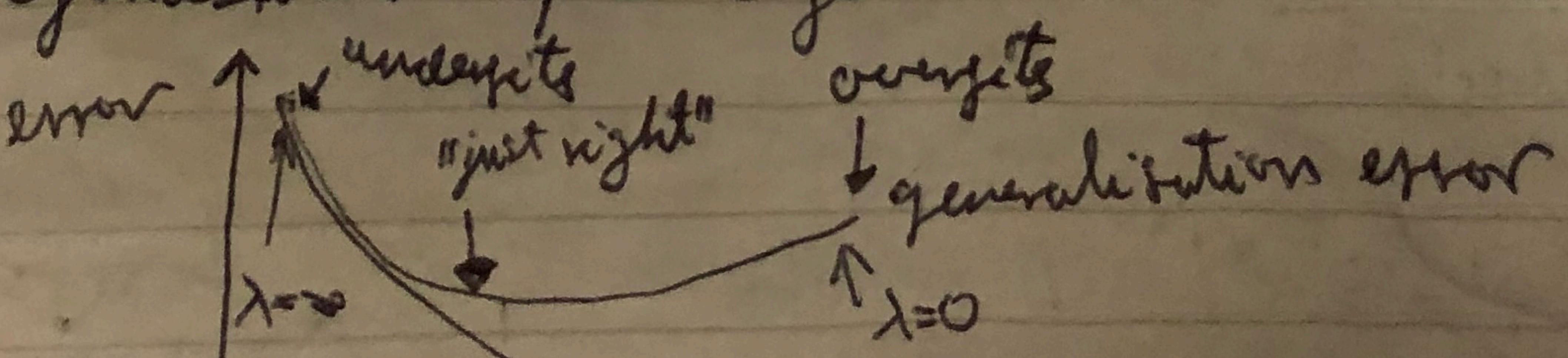
eg log-reg. model.

$$\rightarrow P_g(\Theta) : \Theta \sim \mathcal{N}(0, \tau^2 \mathbb{I}) \Leftrightarrow P_g(\Theta) = \frac{1}{\sqrt{2\pi}^n |\tau^2 \mathbb{I}|^{1/2}} \exp(-\Theta^T (\tau^2 \mathbb{I})^{-1} \Theta)$$

As it turns out: Now, if you plug in  $P_g(\Theta)$  for  $P$  in the above optimisation problem and solve it, our regularisation technique falls out.

Frequentist:  $\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} P(S|\Theta) - \text{MLE}$ : non-regularised per default (but one can apply it after

Bayesian: Prior distribution  $P_g(\Theta) \rightarrow \underset{\Theta}{\operatorname{argmax}} P(\Theta|S) - \text{MAP estimator}$   
- regularised per definition.



model complexity  
(lg degree of polynomial)

## Train/dev/test sets

10,000 examples: Eg  $\Theta_0 + \Theta_1 x$  vs  $\Theta_0 + \Theta_1 x + \Theta_2 x^2$ , etc., or choose  $\epsilon$  or  $C$

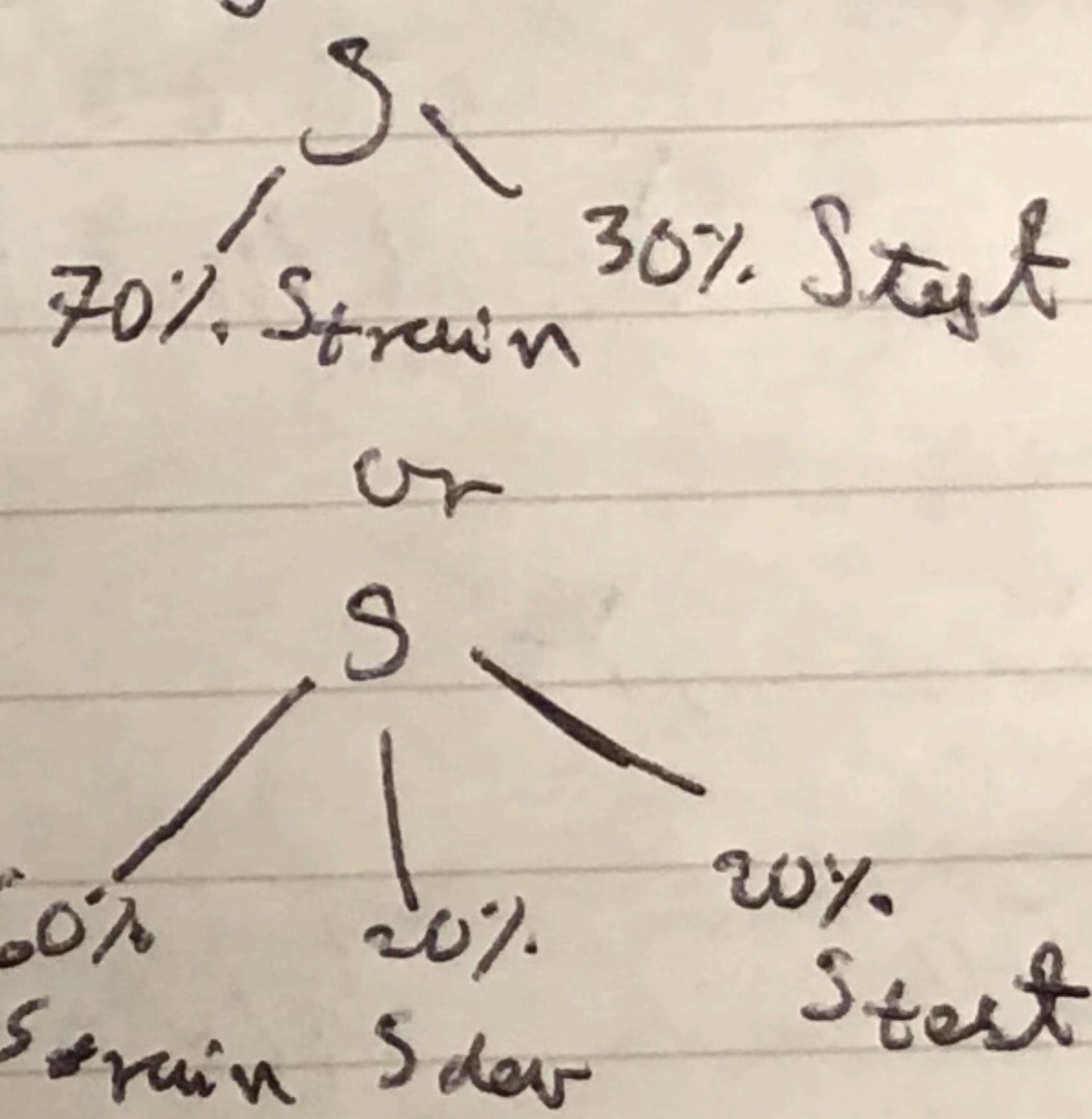
$$S \rightarrow S_{\text{train}}, S_{\text{dev}}, S_{\text{test}}$$

- Train each model (option for degree of polynomial) on  $S_{\text{train}}$ . Get some hypothesis  $h_i$ .
- Measure error on  $S_{\text{dev}}$ . Pick model with lowest error on  $S_{\text{dev}}$ .  
 $\Rightarrow S_{\text{dev}}$  is there to select hyperparameters.
- Optional: Evaluate algorithm on test set  $S_{\text{test}}$ , and report error.

Note: some	degree	$S_{\text{dev}}$ error	$S_{\text{test}}$
tiny, people	1	10	10
full Sdev	2	5.7	5.0
= sum. rd.	3	5.0	5.0
	4	4.9	5.0
	5	7	7
	6	10	10
	:	:	:

### Holdout cross-validation (simple)

Rules of thumb for "reasonably sized datasets":



But for enormous datasets, need smaller % on  $S_{\text{dev}}$ ,  $S_{\text{test}}$ . In general, just pick  $S_{\text{dev}}/S_{\text{test}}$  large enough to differentiate your various models reliably, but not much larger than that.

Note: never base any decisions on  $S_{\text{test}}$ , to remain unbiased.

Small datasets: For 70-30 split: do we really want to use "just" 70 training examples?

Not usually, hence K-fold CV:

$\left| \begin{array}{c} s \\ x^1, y^1 \\ \vdots \\ x^{10}, y^{10} \end{array} \right|$   $k=5$  only for illustration,  $k=10$  is typical.  
 For  $i=1, \dots, k$  train (fit parameters) on  $k-1$  pieces  
 and test on remaining 1 piece, and then average all these.

Optional Final step: Regit model on 100% of data.

Advantage: Uses all data. Disadvantage: Computational expensive

$m=20-100$ :  
 ~~$m=20$~~  a extreme version of k-CV is "leave-out-one-CV":  
 $k=m$ .

### Feature Selection

Start with  $f = \emptyset$

Repeat:

- 1) Try adding each feature  $i$  to  $f$ , and see which single feature addition most improves dev set performance.
- 2) Add that feature to  $f$

$x_1, \dots, x_5$

$$f = \emptyset \quad h_f(x) = \Theta_0$$

$$\begin{cases} \emptyset + x_1 \\ \emptyset + x_2 \\ \vdots \\ \emptyset + x_5 \end{cases}$$

$$h_{\emptyset}(x) = \Theta_0 + \Theta_1 x_5$$

$$\Rightarrow f = \{x_2\}$$

$h_{\emptyset}(x) = \Theta_0 + \Theta_1 x_2$  new model. Trying out  $x_4$   
 $x_4$  adds best performance  $\Rightarrow$  new model is  $\Theta_0 + \Theta_1^T \begin{bmatrix} x_2 \\ x_4 \end{bmatrix}$ .  
 This is a reasonable, greedy algorithm.