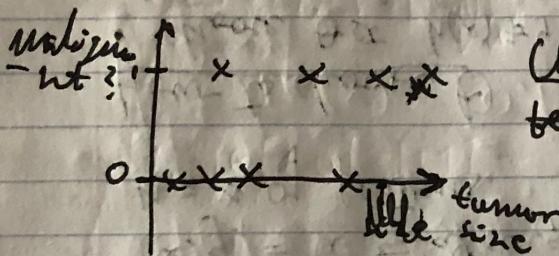
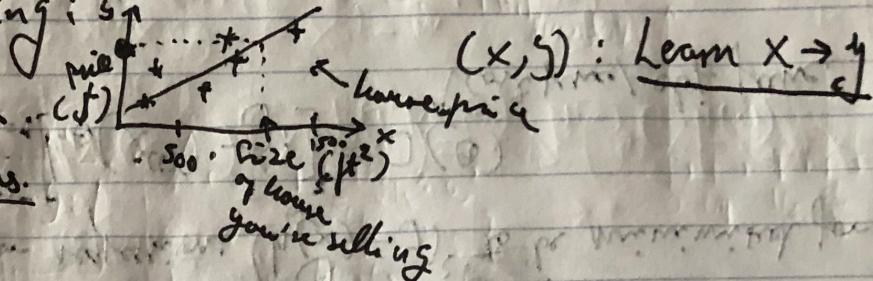


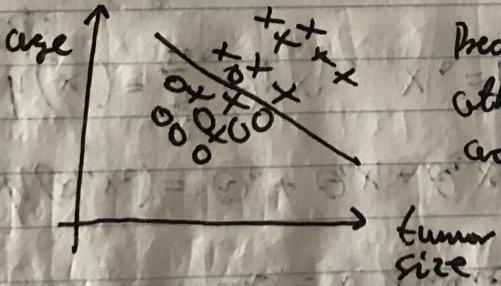
# CS229 Lecture 1

- First significant ML result: was A. Samuel's 1959 Checkers AI.
- Well-posed learning problem: for problem A, computer learns from exp E w.r.t task T acc. to pay. measure P, if its perf. on T, as measured by P, improves with E.
- Supervised learning:  $y = f(x)$

Regression problem: because  $y$  is continuous.



Classification problem because  $y$  is discrete



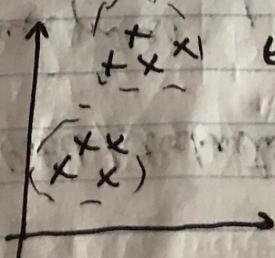
Predict whether tumor malign or benign  
other features: clump thickness, unif. of cell size, adhesion etc.

Support vector machine: by mean of features,  $\mathbf{x}$ -dim vector to rep. a patient ...

Future subjects:

- ML Strategy (learning theory)
- Deep Learning

Unsupervised learning:  
Basically, finds structures and patterns in data.



unlabelled data  
Goal: figure out interesting structures here.

Reinforcement learning: Reward good behavior, punish bad.

Machine learning: linear regression / gradient descent

## Lecture 2

### Outline

- Linear regression; • Batch/stochastic gradient descent; • Normal Eq.

Def: Take house prices:

How to represent  $h$ ?

Let's say  $h(x) = \Theta_0 + \Theta_1 x$

(can add bedrooms as a feature,  $x_2$ .  $(x_0, x_1, x_2)$ )

Set not

- ation:  $\Rightarrow h(x_1, x_2) = \Theta_0 + \Theta_1 x_1 + \Theta_2 x_2$  "parameters" to be trained

$$\Leftrightarrow h(x) = \sum_{j=0}^2 \Theta_j x_j, x_0 = 1 \Leftrightarrow \Theta = \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \Theta_2 \end{bmatrix}, x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \end{bmatrix} [h_\Theta(x) = h(x)]$$

def:

$m$  = # training examples;  $X$  = "inputs/features";  $y$  = "output/target variable"

$(x, y)$  = a training example;  $(x^{(i)}, y^{(i)})$  = i-th training example

$n$  = # of features

### How to choose $\Theta$

• choose  $\Theta$  st  $h(x) \approx y$  in training example

• In ordinary linear regression:  $(h(x) - y)^2$  to be minimized for  $\Theta$ .

$$\Leftrightarrow \text{minimize}_{\Theta} \frac{1}{2} \sum_{i=0}^m [h_\Theta(x^{(i)}) - y_i]^2 = J(\Theta)$$

Cost func/ obj func

### Gradient descent

• Start with some initial  $\Theta$  (say  $\Theta = 0$ )

• keep changing  $\Theta$  to reduce  $J(\Theta)$

• Multi: learning rate increment by one normal equality

$$\Theta_j := \Theta_j - \alpha \frac{\partial}{\partial \Theta_j} J(\Theta); \quad \alpha: a+1; \quad a=b$$

for  $j = 0, 1, 2, \dots$  (up to number of features)

assignment

for one training example. (can ignore the sum since derivation  
is for a fully linear operation)  $\hat{y} = \theta_0 + \theta_1 x_1$

$$\frac{\partial J}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (\hat{y}_i - y_i) x_i^{(i)} = \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x_i^{(i)} - y_i) x_i^{(i)}$$

$$\Rightarrow \theta_j := \theta_j - \alpha \underbrace{\sum_{i=1}^m (\theta_0 + \theta_1 x_i^{(i)} - y_i) x_i^{(i)}}_{\text{keep in mind } \frac{\partial J}{\partial \theta_j}} \leftarrow \text{when remembering that we're dealing with a sum here, in } \Sigma$$

keep in mind  $\frac{\partial J}{\partial \theta_j}$

that in practice, it's always good to normalize  $\vec{x}$  in some way,  $\vec{y}$ , by dividing it by  $\|\vec{x}\|$  (scaling).

Repeat the descent until convergence, for  $j = 0, 1, 2, \dots, n$   
(normal)

Batch gradient descent: computes gradient exactly per iteration  
not suitable when  $m$  is huge. e.g. hundred million ( $m \approx 10^8$ )

Repeat { ~~for all samples~~ for every  $i$   
For  $i=1$  to  $m$  }  $\theta_j := \theta_j - \alpha (\hat{y}_i - y_i) x_i^{(i)}$  ← data is shuffled

Stochastic gradient descent: updates  $\vec{\theta}$  using

one sample at a time (or a

minibatch). Sample  $\vec{x}_i$  is stochastically chosen



involves random path,  
but on avg. moves towards minimum.

Normal Equation: Gives result immediately (solve linear problem)

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_n} \end{bmatrix}. \text{ Let's say } f \in \mathbb{R}^{2 \times 2} \Rightarrow A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix}, f: \mathbb{R}^{2 \times 2} \rightarrow \mathbb{R}$$

$$f \circ g f(A) = A_{11} + A_{12} \Rightarrow f \left( \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} \right) = 5 + 6 = 11$$

$$\nabla_A f(A) = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \frac{\partial f}{\partial A_{12}} \\ \frac{\partial f}{\partial A_{21}} & \frac{\partial f}{\partial A_{22}} \end{bmatrix} \text{ For the above example, it's be } \nabla_A f(A) = \begin{bmatrix} 1 & 2 \\ 0 & 0 \end{bmatrix}$$

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{bmatrix}$$

$\nabla_{\theta} J(\theta) = \vec{0}$ , gives solution to optimization problem.

If  $A$  is square ( $A \in \mathbb{R}^{n \times n}$ ), then  $\text{tr } A = \sum_i A_{ii}$

interesting note: if  $f(A) = \text{tr } AB \Rightarrow \nabla_A f(A) = B^T$ .  $\text{tr } ABC = \text{tr } CAB \cdot \nabla_A \text{tr } AAC^T$   
 $= CAB + C^TA$ .

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$X\theta = \begin{bmatrix} (x^{(1)})^T \\ (x^{(2)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

$$= \begin{bmatrix} x^{(1)\top} \theta \\ x^{(2)\top} \theta \\ \vdots \\ x^{(m)\top} \theta \end{bmatrix} = \begin{bmatrix} h_{\theta}(x^{(1)}) \\ h_{\theta}(x^{(2)}) \\ \vdots \\ h_{\theta}(x^{(m)}) \end{bmatrix} \text{. Then } \vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$$

$\Rightarrow J(\theta) = \frac{1}{2} (\vec{X}\theta - \vec{y})^T (\vec{X}\theta - \vec{y}) \Leftrightarrow \text{sum of squared errors.}$

$$= \begin{bmatrix} h_{\theta}(x^{(1)}) - y_1 \\ \vdots \\ h_{\theta}(x^{(m)}) - y_m \end{bmatrix}$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \frac{1}{2} (\vec{X}\theta - \vec{y})^T (\vec{X}\theta - \vec{y}) = \frac{1}{2} \nabla_{\theta} (\theta^T \vec{X}^T - \vec{y}^T)(\vec{X}\theta - \vec{y})$$

$$= \frac{1}{2} \nabla_{\theta} [\theta^T \vec{X}^T \vec{X} \theta - \theta^T \vec{X}^T \vec{y} - \vec{y}^T \vec{X} \theta + \vec{y}^T \vec{y}]$$

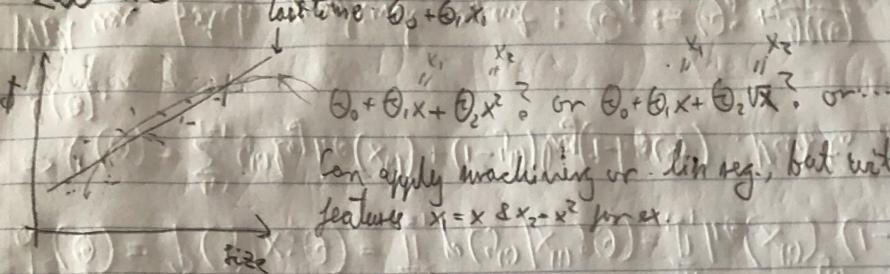
$$\left[ \text{e.g. } (a+b)(a-b) = a^2 - ab - ba + b^2 \right]$$

$$= \frac{1}{2} [\vec{X}^T \vec{X} \theta + \vec{X}^T \vec{X} \theta - \vec{X}^T \vec{y} - \vec{X}^T \vec{y}] = \vec{X}^T \vec{X} \theta - \vec{X}^T \vec{y} \stackrel{\text{set } \vec{0}}{=} \vec{0}$$

$$\vec{X}^T \vec{X} \theta = \vec{X}^T \vec{y} \text{ "normal equations", } \Rightarrow \theta_{\text{opt.}} = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{y}$$

Note: If  $\vec{X}^T \vec{X}$  is singular, you probably have dependent features. In which case you could simply proceed with a pseudoinverse, best is to go into the features and figure out what's wrong.

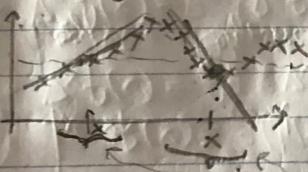
## Lecture 3



## Locally Weighted Regression

"Parametric" learning algo: Fit fixed set of parameters ( $\theta_i$ ) to data. Eg linear.

"Non-parametric" learning algo: Amount of data/parameters you need to keep grows (linearly) with size of data.



To evaluate  $h$  at certain  $x$ :

$$\text{LR: Fix } \theta \text{ to minimize } J(\theta) = \frac{1}{2} \sum_i (y^{(i)} - \theta^T x^{(i)})^2$$

& return  $\theta^T x$ .

$$\text{LWR: Fit } \theta \text{ to minimize } \sum_i w^{(i)} (y^{(i)} - \theta^T x^{(i)})^2$$

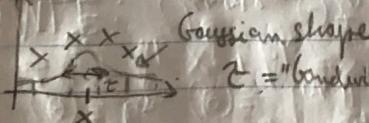
where  $w^{(i)}$  is a "weighting function"

$$w^{(i)} = \exp\left[-\frac{(x^{(i)} - x)^2}{2\sigma^2}\right] \quad \text{tells you how much attention to pay to } y^{(i)}, x^{(i)}$$

Basically, LWR prioritizes accuracy close to the  $x$  where you want to make your prediction.

If  $|x^{(i)} - x|$  is small,  $w^{(i)} \approx 1$ .

If  $|x^{(i)} - x|$  is large,



$z = \text{"bandwidth"}$  has an effect on overfitting

"Not great at extrapolation"

Why least Squares?

Assume  $y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$  ← Error, aka Unmodelled effects/random noise

and  $\epsilon^{(i)} \sim N(0, \sigma^2)$ , IID (big assumption),  $\Rightarrow P(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{\epsilon^{(i)2}}{2\sigma^2}\right]$

$$\Rightarrow P(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}} \Leftrightarrow (y^{(i)} | x^{(i)}; \theta) \sim N(\theta^T x^{(i)}, \sigma^2)$$

least squares falls naturally out of these assumptions.

"likelihood w.r.t.  $\Theta$ "

$$l(\Theta) = p(y|x; \Theta) = \prod_{i=1}^m p(y^{(i)}|x^{(i)}; \Theta) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y^{(i)} - \Theta^T x^{(i)})^2}{2\sigma^2}\right)$$

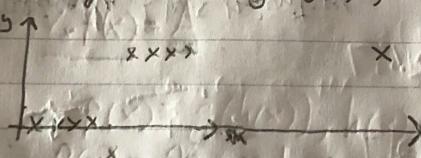
L's Dual interpretation: fixed data = likelihood of  $\Theta$ ; fixed  $\Theta$  = probability of data

"log-likelihood"  
"loss"

$$l(\Theta) = \text{LogL}(\Theta) = \log \prod_{i=1}^m \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y^{(i)} - \Theta^T x^{(i)})^2}{2\sigma^2}\right) = \ln \log \frac{1}{\sqrt{2\pi\sigma^2}} + \sum_{i=1}^m -\frac{(y^{(i)} - \Theta^T x^{(i)})^2}{2\sigma^2}$$

MLE: Choose  $\Theta$  to maximise  $\mathcal{L}(\Theta) / \text{LogL}(\Theta)$ , i.e. choose  $\Theta$  to minimise  $\sum_{i=1}^m (y^{(i)} - \Theta^T x^{(i)})^2 = J(\Theta)$   $\Rightarrow$  choosing  $\Theta$  to minimise  $J$  is equivalent to maximizing probability of data, given iid + N( $\mu, \sigma^2$ ) assumptions.

Classification:  $y \in \{0, 1\}$  (binary classification)

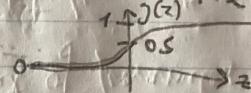


Linear regression is not a good classification algorithm.

Logistic Regression: most commonly used algorithm.

Want  $h_\Theta(x) \in [0, 1]$ .  $h_\Theta(x) = g(\Theta^T x) = \frac{1}{1+e^{-\Theta^T x}}$ ,  $g(z) = \frac{1}{1+e^{-z}}$

= "sigmoid" or "logistic function".



Assume  $P(y=1|x; \Theta) = h_\Theta(x)$   $\Rightarrow P(y=0|x; \Theta) = 1 - h_\Theta(x)$

$$l(\Theta) = P(\vec{y}|x; \Theta) = \prod_{i=1}^m p(y^{(i)}|x^{(i)}; \Theta) = \prod_{i=1}^m h_\Theta(x^{(i)})^{y^{(i)}} (1 - h_\Theta(x^{(i)}))^{1-y^{(i)}}$$

$$\Rightarrow l(\Theta) = \sum_{i=1}^m y^{(i)} \log h_\Theta(x^{(i)}) + (1-y^{(i)}) \log(1-h_\Theta(x^{(i)}))$$

. Need to find  $\Theta$  that maximizes  $l(\Theta)$ .

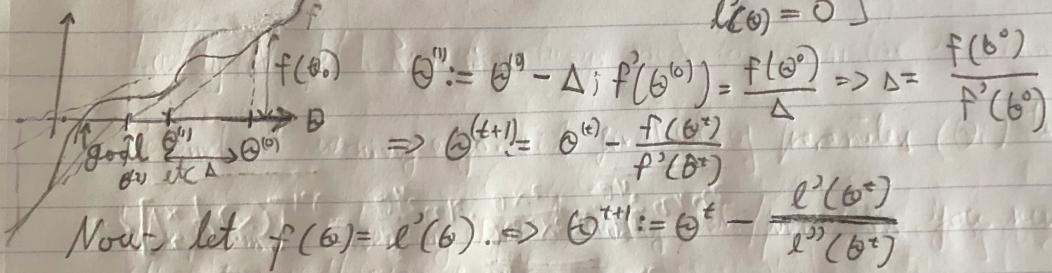
Will use batch gradient ascent:  $\Theta_j := \Theta_j + \frac{\partial l(\Theta)}{\partial \Theta_j}$

$$\Rightarrow \hat{\theta}_j := \theta_j + \alpha \left[ \sum_{i=1}^m (f^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)} \right] = \frac{\partial J}{\partial \theta_j}$$

(PS: logistic regression w/ the sigmoid is concave optimization! Very nice!)

### Newton's Method:

Have  $f$ . Want to find  $\theta$  s.t.  $f(\theta) = 0$ . [or want to maximize  $\lambda(\theta)$ , i.e. want  $\lambda'(\theta) = 0$ ]



$$\theta^{(t)} := \theta^{(t)} - \Delta; f'(\theta^{(t)}) = \frac{f(\theta^{(t)})}{\Delta} \Rightarrow \Delta = \frac{f(\theta^{(t)})}{f'(\theta^{(t)})}$$

$$\Rightarrow \theta^{(t+1)} = \theta^{(t)} - \frac{f(\theta^{(t)})}{f'(\theta^{(t)})}$$

Now, let  $f(\theta) = \lambda'(\theta) \Leftrightarrow \theta^{(t+1)} = \theta^{(t)} - \frac{\lambda''(\theta^{(t)})}{\lambda'''(\theta^{(t)})}$

Newton's method is fast. "Quadratic convergence"  $\Rightarrow$  error  $\propto \varepsilon^2$  for smooth functions

When  $\theta$  is a vector:  $\theta^{(t+1)} = \theta^{(t)} + H^{-1} \nabla \lambda$ , where  $H$  is the Hessian.

$$H_{ij} \in \frac{\partial^2 \lambda}{\partial \theta_i \partial \theta_j} \leftarrow \text{expensive to compute for many vars.}$$

$\Rightarrow$  If  $\lambda$  goes not too much, use Newton's method. Otherwise, <sup>15</sup> gradient descent.  
Or just use quasinewton...

## [Very messy temp-lecture]

### Lecture 4

Topics: Perceptron, Exponential family, Generalised Linear Models, Sigmoid Regression (Multi-class classification).

#### Perceptrons

Logistic Regression: sigmoid function

Perceptron

$$g_p(z) = \begin{cases} 1, & z \geq 0 \\ 0, & z < 0 \end{cases}, h(x) = g_p(\theta^T x)$$

$$g(z) = \frac{1}{1+e^{-z}}$$

a scalar

gradient of relevant CF?

$$\theta_j := \theta_j + \alpha (y_i - h_\theta(x_i)) x_{ij} \leftarrow \text{Perception update?}$$

Not used

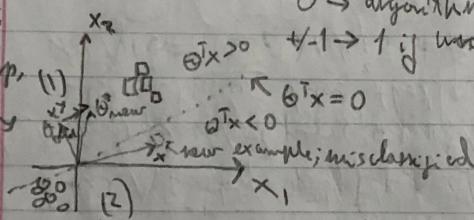
in practice:

high prob-

abil interp.

and other

limitations



$\rightarrow$  algorithm got it right

$\theta^T x > 0 \rightarrow +1$  if wrong  $y=1$ ;  $-1$  if wrong  $y=0$ .

(actually,  $\theta^T x$ )

We want:  $\theta^T x / y = 1$  and  $\theta^T x / y = 0$

$\Rightarrow$  Just move  $\theta$  towards  $x$  when  $y=1$ , or move away when  $y=0$ .

### Exponential families

Eg.: PDF can be written in form:  $p(y; \eta) = b(y) \exp[\eta^T T(y) - a(\eta)]$

$y$  - Data.  $\eta$  - natural parameter.  $T(y)$  - sufficient statistic  $\equiv y$  for today.

$b(y)$  - base measure;  $a(\eta)$  - log-position (log, because  $p(y; \eta) = \frac{b(y) \exp[\eta^T T(y)]}{a(\eta)}$ )

Ex: Bernoulli (Binary data):  $\phi = \text{probability of event. Is it part of expon. families?}$

$$p(y; \eta) = \phi^y (1-\phi)^{1-y} = \exp\left(\log(\phi) (1-\phi)^{-1}\right) = \exp\left[\frac{\log \phi}{1-\phi} y + \log(1-\phi)\right].$$

From this, we see that:  $b(y)=1$  &  $T(y)=y$  &  $\eta = \log \frac{\phi}{1-\phi} \Leftrightarrow \phi = \frac{e^\eta}{1+e^\eta}$  &  $a(\eta) = \log(1+e^\eta)$

Ex: Gaussian (w/ fixed covariance):

$$\text{Let } \sigma^2 = 1. p(y; \mu) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y-\mu)^2}{2}\right) = \frac{1}{\sqrt{2\pi}} e^{-\frac{y^2-2\mu y + \mu^2}{2}} \exp\left[\frac{\mu y - \mu^2/2}{\sqrt{2\pi} T(y)}\right]$$

$$b(y) = \frac{1}{\sqrt{2\pi}} \exp(-y^2/2); T(y) = y, \mu = \mu, a(y) = \mu^2/2$$

when to use  
this model?  
I think in

Real value data: Gaussian  
Binary-valued: Bernoulli  
Count: Poisson  
R<sup>+</sup>: Gamma/exponential

Data over distri: Beta, Dirichlet, Bayesian

why we like exponentials

## Properties

(a) MLE wrt  $\gamma$  is concave / NLL is convex

b)  $E[y|\gamma] = \partial a(\gamma)/\partial \gamma$  no entry

c)  $\text{Var}[y|\gamma] = a(\gamma) \gamma^2$  - rotation! :)

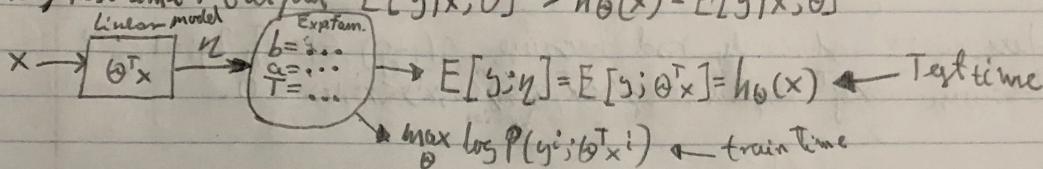
## Generalised Linear models (GLM)

Assumptions/design choices:

i)  $y|x; \theta \sim \text{Exponential Family}(\eta)$

ii)  $\eta = \theta^T x$ ,  $\theta \in \mathbb{R}^n$ ,  $x \in \mathbb{R}^n$

iii) Test time: output  $E[y|x; \theta] \Rightarrow h_\theta(x) = E[y|x; \theta]$



## GLM Training

Learning update rule:  $\theta_j := \theta_j + \alpha(y^i - h_\theta(x^i))x_j$

this apparently generically  
falls out when taking  
gradients etc.

Terminology  $\eta$  - "natural parameter";  $E(y|\gamma) = g(\eta) \rightarrow$  "canonical response func."

$\eta = g^{-1}[E(y|\gamma)]$  - "link function". [also, from prop. b:  $g(\eta) = \partial a/\partial \eta$ ]

## Parametric names

$\theta$  - model param.,  $\eta$  - natural param. Canonical params:  $\begin{cases} \text{Binomial} \\ \text{Gauss} \\ \text{Poisson} \end{cases}$

Learn -  $\theta^T x \rightarrow$   
Design choice

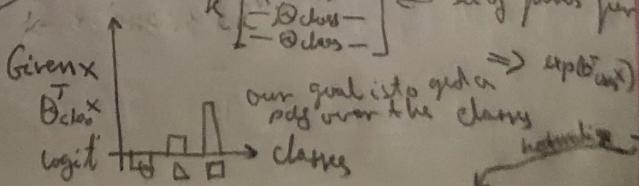
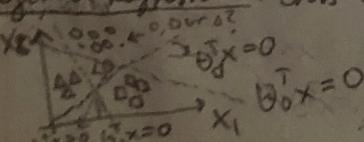
Logistic regression  $h_\theta(x) = F(y|x; \theta) = \phi = \frac{1}{1+e^{-\eta}}$  sigmoid pops out internally from this machinery

Note: Task decides your choice of distribution.

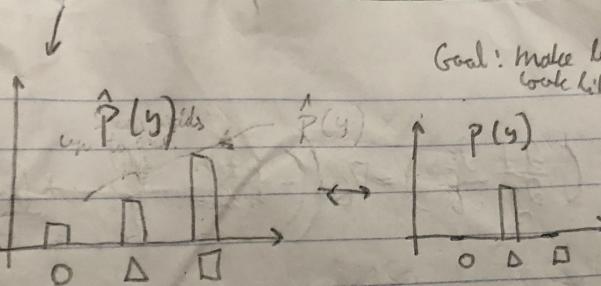
Assumptions  $\Rightarrow$  baseline this job is sampled from the Gaussian Regression  $x \leftarrow$  find  $\hat{x}$  etc...  $\theta^T x = y$

$k \rightarrow$  classes;  $x \in \mathbb{R}^n$ ; label  $y = \sum \beta_k \cdot g[0, 0, 1, 0] \leftarrow$   
 $\theta$  class  $\in \mathbb{R}^n$ ;  $K$  each class  $\in \{0, 1, 2, \dots\}$   
 $k \left[ \begin{array}{c} \text{class } 1 \\ \text{class } 2 \\ \vdots \\ \text{class } k \end{array} \right] \leftarrow$  set of params per class

Softmax Regression: cross-entropy



$$\sum_{i \in \{0, \Delta, \square\}} e^{\theta_i^T x}$$



Goal: make left dist  
look like right

$$\text{Cross Ent}(\mathbf{p}, \hat{\mathbf{p}}) = \sum_{y \in \{0, \Delta, \square\}} p(y) \log \hat{p}(y) = -\log \hat{p}(y_0) = -\log \frac{e^{\theta_0^T x}}{\sum_{i \in \{0, \Delta, \square\}} e^{\theta_i^T x}}$$

treat as  
loss & do gradient  
descent.