

## Grading Scheme

Quizzes 20%, Assgn 30%, MidSem 20%, EndSem 30%  
 (24<sup>th</sup> Aug, 19<sup>th</sup> Oct)

Quizzes → Short quizzes (~2, 10% + 10%)

Assignments → ~3 A, 2-3 weeks, Groups of 5

★ Machine Learning is the art and science of solving ill-understood tasks.

→ Well-understood tasks → { RJ 23 PA 2262  
 Rajasthan { RJ 23 SM 2824  
 Uttar Pradesh { UP 67 AB 3601  
 Delhi { DL 7C MD 750

→ Ill-understood

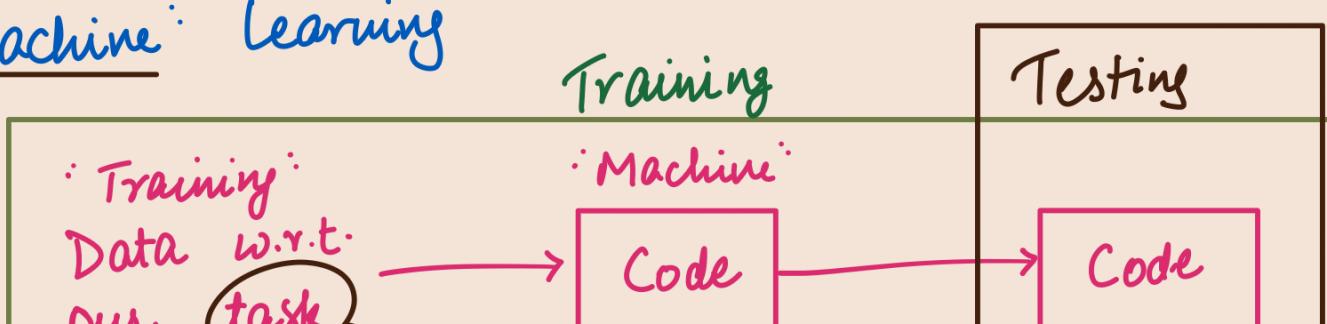
R R → R      D D D D → D

Difficult to formally explain the rules of classification.

Moreover, humans can mostly do only those tasks which they are already familiar with.

Why is it called

Machine Learning

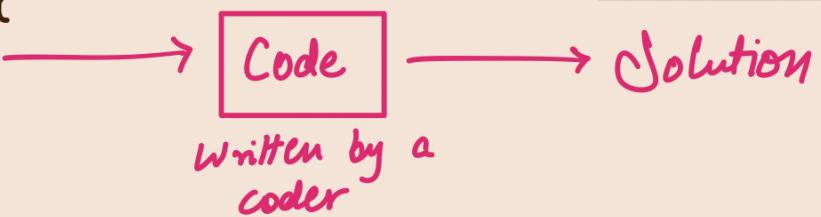


our work  
→  
Ill understood

Algorithm  
written by a  
human

Model  
produced  
automatically  
by the ML algo

Well understood  
task



- \* Nature is governed by laws. Humans are sometimes able to discover these with a lot of time and effort.

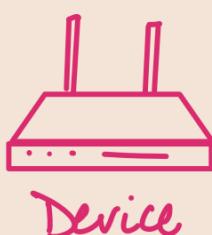
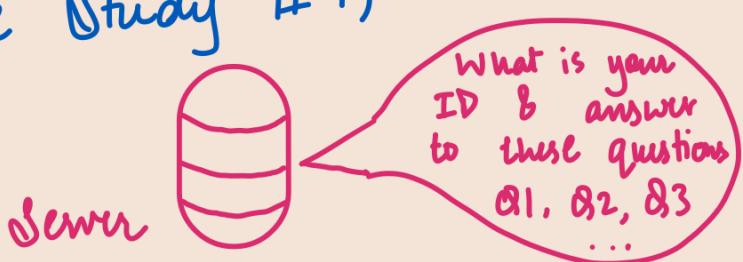
ML works by discovering laws too.

$$P = \boxed{\dots} \approx \boxed{\dots} + \boxed{\dots}$$

ML works by analyzing data to identify laws or patterns that seem to explain the data well.

The laws learnt by ML may be too complex to be interpreted by humans.

- Authentication using secret questions  
(Case Study #1)

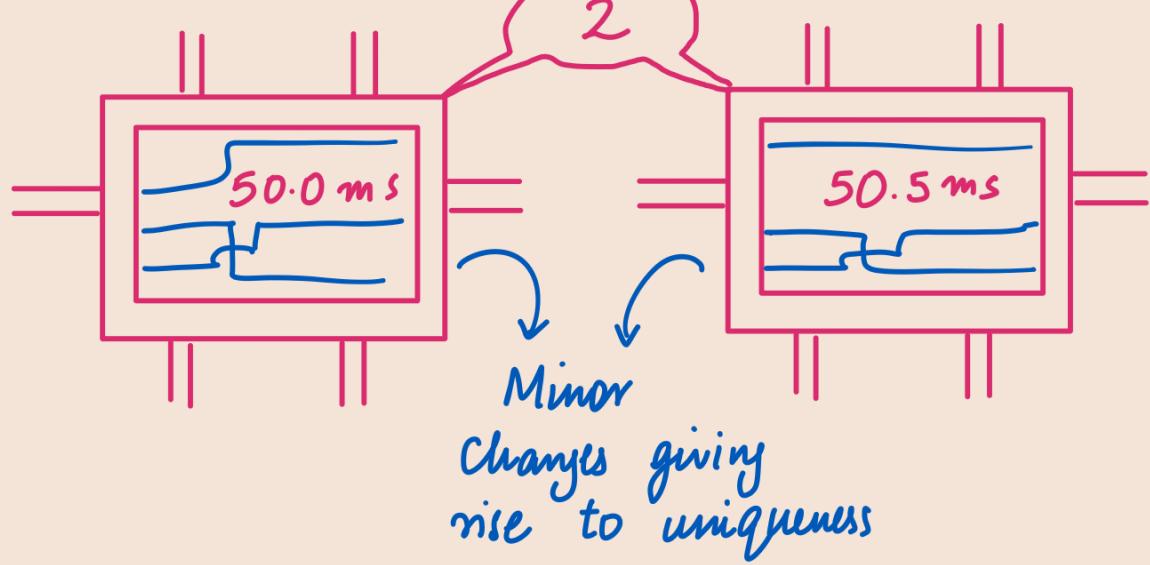


Questions → Should be unique w.r.t. Device

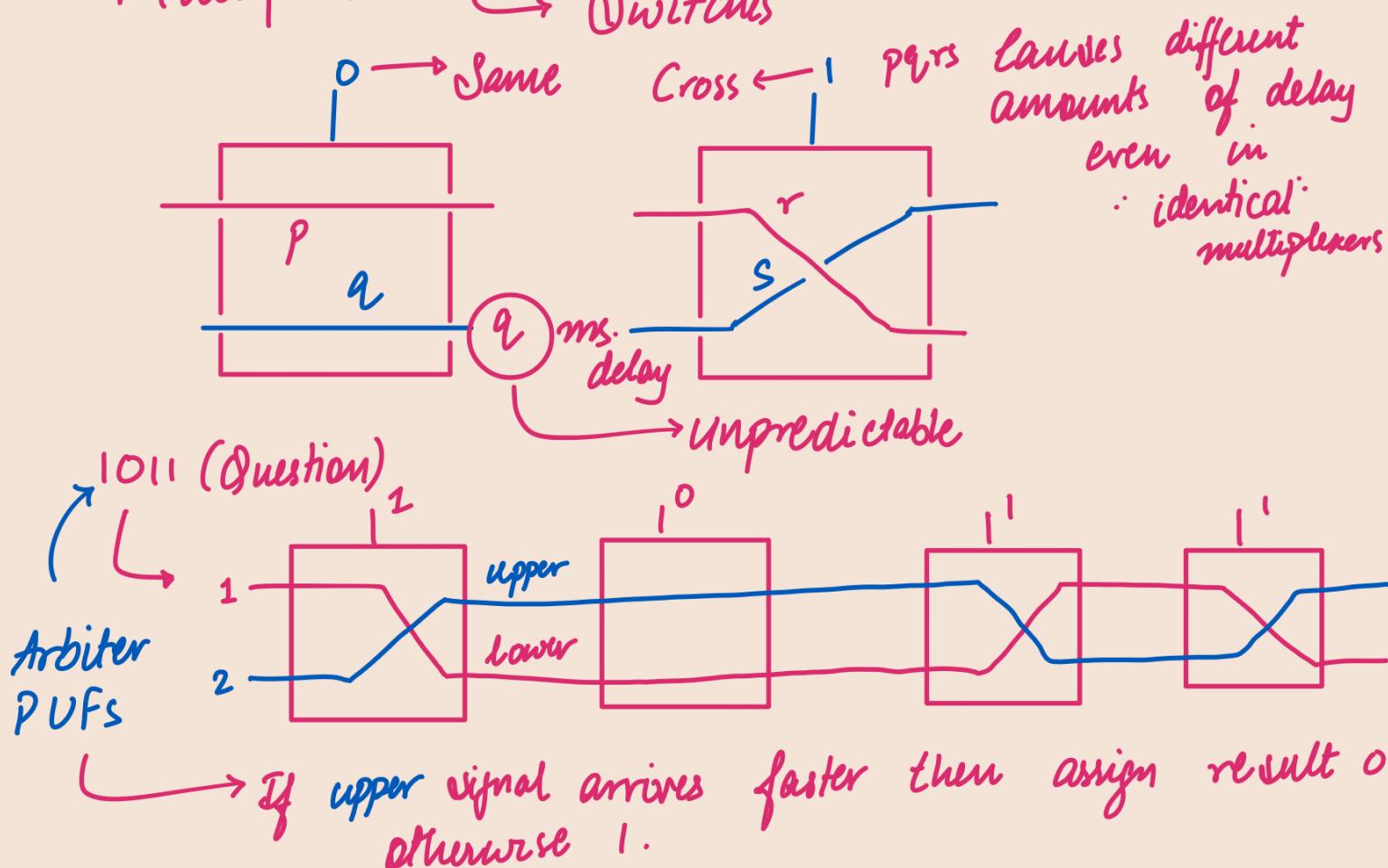
PUFs

Physically Unclonable Functions

$$1+1=?$$



\* Multiplexer  $\rightsquigarrow$  Switches



- Questions — Challenges ; Answers — Responses
- Multiple Multiplexers  $\rightarrow$  Circuits more complex  
 $\Rightarrow$  Difficult to predict
- Also Yes the number of questions  
Usually 64 multiplexers  
 $\rightarrow 2^{64}$  questions
- Knowing the answer to a few questions

$\Rightarrow$  Knowing / Ability to predict answers to other questions: WRONG

An attacker can see responses on a few challenges & use ML to predict responses on all other challenges. No matter if # bits = 32/64

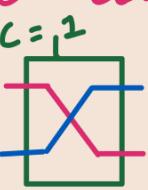
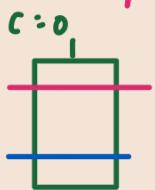
Assume 64 bits of 0...63 labels}

★  $t_i^u \rightarrow$  Time taken for upper signal to reach out in the  $i^{\text{th}}$  ( $i \in \{0, 1, \dots, 63\}$ ) multiplexer  
 $t_i^l \rightarrow \dots$  lower signal ...

Observe that answer is 0 if  $t_{63}^u < t_{63}^l$ , 1 otherwise  
 Also, note that  $t_i^u$  and  $t_i^l$  depend on  $t_0^u$ ,  $t_0^l$ ,  $p_i$ ,  $q_i$ ,  $r_i$ ,  $s_i$  and  $c_i$ .

$c_i$  dictates if  $t_0^l/t_0^u$  (previous delay) will get carried forward in which branch &  $p_i, q_i, r_i, s_i$  give us the delay introduced by the  $i^{\text{th}}$  mux itself.

More concretely,



$$t_i^u = (1 - c_i) \cdot (t_0^u + p_i) + c_i \cdot (t_0^l + s_i)$$

$$t_i^l = (1 - c_i) \cdot (t_0^l + q_i) + c_i \cdot (t_0^u + r_i)$$

Let  $\Delta_i = t_i^u - t_i^l$  denote the lag  
 All that matters is whether the upper signal reaches first or not  $\Rightarrow \Delta_{63} < 0$  or not

$$\alpha_i := (p_i - q_i + r_i - s_i)/2, \beta_i := (p_i - q_i - r_i + s_i)/2,$$

$$\Delta_{-1} = 0, \quad d_i := 1 - 2c_i, \quad d_i \in \{-1, 1\}$$

$$\Delta_i = d_i \cdot \Delta_{i-1} + \alpha_i \cdot d_i + \beta_i$$

$$\Delta_0 = \alpha_0 \cdot d_0 + \beta_0$$

$$\Delta_1 = \Delta_0 \cdot d_1 + \alpha_1 \cdot d_1 + \beta_1$$

$$= \alpha_0 \cdot d_1 \cdot d_0 + (\alpha_1 + \beta_0) \cdot d_1 + \beta_1$$

$$\Delta_2 = \Delta_1 \cdot d_2 + \alpha_2 \cdot d_2 + \beta_2$$

$$= \alpha_0 \cdot d_2 \cdot d_1 \cdot d_0 + (\alpha_1 + \beta_0) \cdot d_2 \cdot d_1 + (\alpha_2 + \beta_1) \cdot d_2 + \beta_2$$

$$\dots \Rightarrow \Delta_{63} = w_0 \cdot x_0 + w_1 \cdot x_1 + \dots + w_{63} \cdot x_{63} + \beta_{63}$$

$$\Delta_{63} = w^T x + \beta$$

where

$$\begin{cases} x_i = d_i d_{i+1} \dots d_{63} \\ w_0 = \alpha_0 \\ w_i = \alpha_i + \beta_{i-1} \quad (i > 0) \end{cases}$$

$\Delta_{63} < 0 \Rightarrow$  upper signal wins  
 $\Rightarrow$  answer = 0

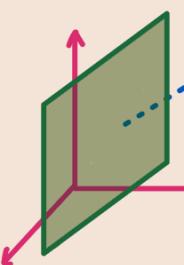
$\Delta_{63} > 0 \Rightarrow t_{63}^u > t_{63}^l \Rightarrow$  lower wins  $\Rightarrow$  answer = 1

$\Rightarrow$  Answer is simply  $\frac{\text{sign}(w^T x + \beta) + 1}{2}$

Nothing but a **linear classifier**  $\rightarrow$  single vector (d-dim.)  
 $\hookrightarrow$  Model  $\equiv \{w, b\}$   $\rightarrow$  scalar term (bias)  
features are too

$\hookrightarrow$  Predict  $x$  by checking if  $w^T x + b > 0$

Decision boundary: hyperplane ( $w^T x + b = 0$ )



$w$ :  $\perp$  / normal vector to the hyperplane

Let  $x, y$  be vectors on hyperplane then

$$w^T x + b = 0 = w^T y + b \Rightarrow w^T (x - y) = 0$$

$\Rightarrow x - y \perp$  to  $w$  &  $\parallel$  to the hyperplane

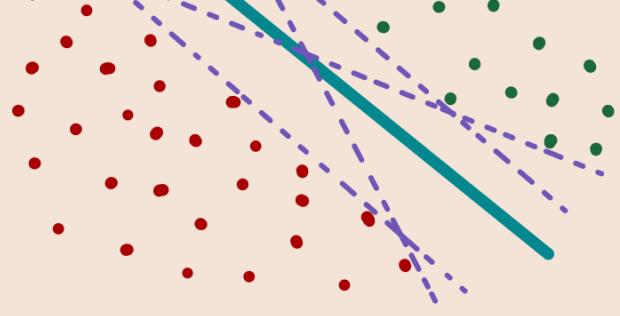
$b$  shifts the plane } can be thought of as a threshold that  $w^T x$  has to pass to make the decision 1 }

Best: Linear Classifier

Better not to select



a model whose decision boundary passes very close to a training data point



Fact: Dist. of origin from hyperplane  $w^T x + b = 0$  is  $|b| / \|w\|_2$

$$\|v\|_p = \left( \sum |v_i|^p \right)^{1/p} \Rightarrow \|w\|_2 = \text{length of } w$$

Fact: Dist. of a point  $p$  from " is  $|w^T p + b| / \|w\|_2$

Given train data  $\{(x^i, y^i)\}_{i=1}^n$  for a binary classification problem where  $x^i \in \mathbb{R}^d$  and  $y^i \in \{-1, 1\}$ , we want two things from the classifier:

Demand 1: Classify every point correctly:

$$\rightarrow \forall i \in [n] \quad \text{sign}(w^T x^i + b) = y^i$$

$$\rightarrow \forall i \in [n] \quad y^i \cdot (w^T x^i + b) \geq 0$$

Demand 2: Not let any data point come close to the boundary

$\rightarrow \min_{i \in [n]} |w^T x^i + b| / \|w\|_2$  be as large as possible

## Support Vector Machine

(S.V.M.)

Please find me a linear classifier that {satisfies D1 & D2}.

Mathematical way  $\rightarrow$

$$\max_{w, b} \left\{ \min_{i \in [n]} \left\{ |w^T x^i + b| / \|w\|_2 \right\} \right\}$$

such that  $y^i \cdot (w^T x^i + b) \geq 0 \quad \forall i \in [n]$

as an Optimization Problem w/ an objective and lots of constraints

Set of all points satisfying all constraints is

called the "feasible" set of the optimization problem

Now, assume  $\exists$  params that perfectly classify all train data. Consider one such  $(w, b)$ :

$$|w^T x^i + b| / \|w\|_2 = \frac{|y_i \cdot (w^T x^i + b)|}{\|w\|_2} \quad \left\{ \because y_i \in \{-1, 1\} \right.$$

$\min_{i \in [n]} \frac{|y_i \cdot (w^T x^i + b)|}{\|w\|_2} \equiv \min_{i \in [n]} \frac{y_i \cdot (w^T x^i + b)}{\|w\|_2}$

Model has perfect classification

Will remove the assumption that train data is linearly separable. Ex: ::

Let  $i_0$  be the data point that comes closest:

$$\Rightarrow \min_{i \in [n]} \frac{y^i \cdot (w^T x^i + b)}{\|w\|_2} = \frac{y^{i_0} \cdot (w^T x^{i_0} + b)}{\|w\|_2} \quad \left\{ \begin{array}{l} \text{Assuming} \\ (w, b) \text{ is a} \\ \text{perfect classifier} \end{array} \right.$$

Let  $\epsilon = y^{i_0} \cdot (w^T x^{i_0} + b)$

consider  $\tilde{w} = w / \epsilon$ ,  $\tilde{b} = b / \epsilon$

$$\Rightarrow y^i \cdot (\tilde{w}^T x^i + \tilde{b}) \geq 1 \quad \forall i \in [n]$$

$\min_{i \in [n]} \frac{y^i \cdot (\tilde{w}^T x^i + \tilde{b})}{\|\tilde{w}\|_2} = \frac{1}{\|\tilde{w}\|_2} \quad \left\{ \because y^{i_0} \cdot (\tilde{w}^T x^{i_0} + \tilde{b}) = 1 \right.$

$\Rightarrow$  Thus, instead of searching for  $(w, b)$ , easier to search for  $(\tilde{w}, \tilde{b})$

$$\max_{\tilde{w}, \tilde{b}} \frac{1}{\|\tilde{w}\|_2} \rightarrow \min_{\tilde{w}, \tilde{b}} \|\tilde{w}\|_2 \rightarrow \min_{\tilde{w}, \tilde{b}} \|\tilde{w}\|_2^2 \rightarrow \min_{\tilde{w}, \tilde{b}} \frac{\|\tilde{w}\|_2^2}{2}$$

s.t.  $y^i (\tilde{w}^T x^i + \tilde{b}) \geq 1 \quad \forall i \in [n]$

If a linear classifier cannot perfectly classify

data, then,

## C-SVM

$$\min_{\tilde{w}, \tilde{b}, \{\xi_i\}} \frac{1}{2} \|\tilde{w}\|_2^2 + C \sum_{i=1}^n \xi_i$$

$$\text{s.t. } y^i \cdot (\tilde{w}^T x^i + \tilde{b}) \geq 1 - \xi_i \quad \forall i \in [n]$$

$$\text{as well as } \xi_i \geq 0 \quad \forall i \in [n]$$

Slack variables

Allow some data points to come close to the hyperplane or be misclassified altogether; Having  $\xi_i > 0$  prevents us from artificially inflating the margin from 1.

What if slack variables allow the model to misclassify each data point? C prevents it. If C is set to be 1, then it penalizes the models that slack too much. C is a hyper-parameter

Now,  $\xi_i$  basically allow us to have

$$y^i \cdot (\tilde{w}^T x^i + \tilde{b}) < 1 \quad (\text{even } < 0).$$

No, slack  $\xi_i$  is just  $\xi_i = 1 - y^i \cdot (\tilde{w}^T x^i + \tilde{b})$

But,  $\xi_i \geq 0$  too

∴ If we have  $y^i \cdot (\tilde{w}^T x^i + \tilde{b}) \geq 1$ , then we don't need any slack, i.e.,  $\xi_i = 0$ .

Thus,

$$\xi_i = [1 - y^i \cdot (\tilde{w}^T x^i + \tilde{b})]_+$$

where  $[x]_+ = \max\{x, 0\}$

Suppose on  $(x, y)$ ,  
 $y \in \{-1, +1\}$ , model  
 $(w, b)$  predicts s  
 $\{s = w^T x + b\}$ ,

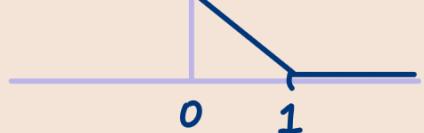
Hinge-Loss  
Function

we want  $s \cdot y \geq 0$  for correct classification  
but we also want  $s \cdot y > 0$  for large margin — hinge loss captures both:

$$l_{\text{hinge}}(s, y) = [1 - s \cdot y]_+$$

$$= \begin{cases} 0 & s \cdot y \geq 1 \\ 1 - s \cdot y & \text{otherwise} \end{cases}$$

$$\begin{cases} 0, & \text{if } s \cdot y \geq 1 \\ 1 - s \cdot y, & \text{if } s \cdot y \leq 1 \end{cases}$$



∴ Replacing  $\epsilon_i$  by hinge, we get,

C-SVM :  $\min_{\tilde{w}, \tilde{b}} \frac{1}{2} \|\tilde{w}\|_2^2 + C \sum_{i=1}^n l_{\text{hinge}}(\tilde{w}^\top x^i + \tilde{b}, y^i)$

**Optimization** problems generally consist of three components:

- 1) **Objective**
  - 2) **Constraints**
  - 3) **Regularizer**
- Sometimes, regularizers are taken inside the objective

## \* Using Calculus for Optimization

Method 1: **First Order Optimality Condition**

- Exploits the fact that gradient must vanish at a local optimum.
- For convex functions, local minima are global  
↑ Only for simple functions w/ no constraints
- $f' = 0$       b       $f'' \geq 0$
- Not very practical for checking stability

M2: **(Sub) gradient Descent**

- Direction opposite to gradient offers steepest descent. Choose a subgradient  $g^t$  & learning rate / step length  $\eta_t$  & update  $w^{t+1} \leftarrow w^t - \eta_t g^t$
- Helpful in extending gradient descent to "not-nice" / non-differentiable functions

Some **Concerns**:

- Might get stuck at local minima / saddle points  
⇒ Difficult to find good initialization conditions

- Step length:  $\uparrow$  overshoot,  $\downarrow$  slow convergence
- How to define and decide convergence?

M 3: Creating a dual problem (constrained optimization)  
We wish to solve:

$$\begin{array}{l} \min_{x \in \mathbb{R}^d} f(x) \\ \text{s.t. } g(x) \leq 0 \end{array} \xrightarrow{\substack{\text{integrate constraints} \\ \text{into optimization} \\ \text{use Barrier (indicator) function}}} \min_{x \in \mathbb{R}^d} f(x) + r(x)$$

Karush-Kuhn-Tucker

$$r(x) := \begin{cases} 0 & : g(x) \leq 0 \\ \infty & : \text{otherwise} \end{cases}$$

$$r(x) = \max_{\alpha \geq 0} \alpha \cdot g(x)$$

$$\min_{x \in \mathbb{R}^d} \left\{ \max_{\alpha \geq 0} \{ f(x) + \alpha \cdot g(x) \} \right\}$$

Some steps for handling multiple + equality constraints:

S1:	$\max f$	$\longrightarrow$	$\min -f$	}	Only want min.
S2:	$g \geq 0$	$\longrightarrow$	$-g \leq 0$		
S3:	$s = 0$	$\longrightarrow$	$s \leq 0$		
S4:	C constraints	$\longrightarrow$	C variables		

dual variables / ←  
Lagrange multipliers → Original variables  
called primal variables

The Lagrangian

$$\begin{array}{l} \min_{x \in \mathbb{R}^d} f(x) \\ \text{s.t. } g_1 \leq 0, g_2 \leq 0 \dots g_c \leq 0 \end{array} \xrightarrow{\quad} \min_{x \in \mathbb{R}^d} \left\{ \max_{\substack{\alpha \in \mathbb{R}^c \\ \alpha_c \geq 0}} \left\{ f(x) + \sum_{c=1}^C \alpha_c \cdot g_c(x) \right\} \right\}$$

$$\mathcal{L}(x, \alpha) = f(x) + \sum_{c=1}^C \alpha_c \cdot g_c(x)$$

$$\max_{\substack{\alpha \in \mathbb{R}^c \\ \alpha_c > 0}} \{ \mathcal{L}(x, \alpha) \} = \begin{cases} \infty, & x \text{ violates even one constraint} \\ f(x), & x \text{ satisfies every single constraint} \end{cases}$$

$$\min_{x \in \mathbb{R}^d} \left\{ \max_{\substack{\alpha \in \mathbb{R}^c \\ \alpha_c > 0}} \left\{ f(x) + \sum_{c=1}^C \alpha_c \cdot g_c(x) \right\} \right\} \quad \text{Primal Problem}$$

↓ If original problem is convex and "nice"

$$\begin{array}{l} \hookrightarrow \text{Solution: } (\hat{x}^P, \hat{\alpha}^P) \\ \hookrightarrow \text{Solution: } (\hat{x}^D, \hat{\alpha}^D) \end{array}$$

$$\max_{\substack{\alpha \in \mathbb{R}^c \\ \alpha_c > 0}} \left\{ \min_{x \in \mathbb{R}^d} \left\{ f(x) + \sum_{c=1}^C \alpha_c \cdot g_c(x) \right\} \right\} \quad \text{Dual Problem}$$

**Strong duality:**  $\hat{x}^P = \hat{x}^D$  if original problem is convex and "nice".

**Complementary Slackness:**  $\hat{\alpha}_c^D \cdot g_c(\hat{x}^D) = 0 \quad \forall \text{ constraints } c$

## Hard SVM (w/o a bias)

$$\min_{w \in \mathbb{R}^d} \frac{1}{2} \|w\|_2^2 \quad \text{s.t.} \quad 1 - y^i \cdot w^T x^i \leq 0 \quad \forall i \in [n]$$

$n$  constraints  $\Rightarrow n$  dual vars.

Lagrangian:  $\mathcal{L}(w, \alpha) = \frac{1}{2} \|w\|_2^2 + \sum_{i \in [n]} \alpha_i (1 - y^i \cdot w^T x^i)$

Primal:  $\underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \underset{\alpha \geq 0}{\operatorname{argmax}} \left\{ \mathcal{L}(w, \alpha) \right\} \right\}$

Dual:  $\underset{\alpha \geq 0}{\operatorname{argmax}} \left\{ \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \mathcal{L}(w, \alpha) \right\} \right\}$

Simplifying

$$\underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \frac{1}{2} \|w\|_2^2 + \sum_{i \in [n]} \alpha_i (1 - y^i \cdot w^T x^i) \right\}$$

Unconstrained problem w/ a convex & diff. objective  
 ↪ Apply first order optimality to solve it  
 $\Rightarrow \nabla = 0 \Rightarrow w = \sum_{i \in [n]} \alpha_i y^i x^i$

⇒ Simplifies to

This is the  $\underset{\alpha \geq 0}{\operatorname{argmax}} \left\{ \sum_{i \in [n]} \alpha_i - \frac{1}{2} \sum_{i \in [n]} \sum_{j \in [n]} \alpha_i \alpha_j y^i y^j \langle x^i, x^j \rangle \right\}$

problem solved (libsvm, sklearn) solve

Once we get  $\alpha_i \rightarrow$  we can get w

## Support Vectors

- We have  $\alpha_i$  for every data point
- After solving the data points for which  $\alpha_i \neq 0$  : Support Vectors

(usually  $\ll n$  in number)

- From complementary slackness, we have:  
A data point can be SV only if  $y^i w^T x^i = 1$
- If we have bias  $b$  as a slack variable, then the dual looks like:

$$\underset{\alpha \in \mathbb{R}^n}{\operatorname{argmax}} \left\{ \sum_{i \in [n]} \alpha_i - \frac{1}{2} \sum_{i \in [n]} \sum_{j \in [n]} \alpha_i \alpha_j y^i y^j \langle x^i, x^j \rangle \right\}$$

$$\text{s.t. } \alpha_i \in [0, C] \quad \text{&} \quad \underbrace{\sum_{i \in [n]} \alpha_i y^i}_{=0}$$

$$w = \sum_{i \in [n]} \alpha_i y^i x^i \quad \left. \begin{array}{l} \text{Force on hyperplane} = 0 \\ \text{Torque on hyperplane} = 0 \end{array} \right\}$$

$\sum \alpha_i y^i = 0$  difficult to solve because all  $\alpha_i$  related  $\rightarrow$  Need Sequential Minimal Optimization (SMO)  
by John Platt

- If we omit bias (hide in model vector), dual is

$$\underset{\alpha \in \mathbb{R}^n}{\operatorname{argmax}} \left\{ \sum_{i \in [n]} \alpha_i - \frac{1}{2} \sum_{i \in [n]} \sum_{j \in [n]} \alpha_i \alpha_j y^i y^j \langle x^i, x^j \rangle \right\} \text{s.t. } \alpha_i \in [0, C]$$

- We can solve SVM (no bias) by solving:

primal:  $\underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \left\{ \frac{1}{2} \|w\|_2^2 + C \cdot \sum_{i \in [n]} [1 - y^i w^T x^i]_+ \right\}$

OR

dual:  $\underset{\alpha \in \mathbb{R}^n}{\operatorname{argmax}} \left\{ \sum_{i \in [n]} \alpha_i - \frac{1}{2} \sum_{i \in [n]} \sum_{j \in [n]} \alpha_i \alpha_j y^i y^j \langle x^i, x^j \rangle \right\}$

s.t.  $\alpha_i \in [0, C]$

Approaches like sub / gradient ascent / descent, coordinate ascent / descent can be used.

Will do coordinate maximization for dual optimization var. is  $\alpha \rightarrow$  Take one co-ordinate, i.e., one data point at each step.

$\Rightarrow$  Concentrate only on the terms involving  $\alpha_i$ , we get,

$$\underset{\alpha_i \in \mathbb{R} \text{ s.t. } \alpha_i \in [0, C]}{\operatorname{argmax}} \left\{ \alpha_i - \frac{1}{2} \alpha_i^2 \|x^i\|_2^2 - \alpha_i y^i \sum_{j \neq i} \alpha_j y^j \langle x^i, x^j \rangle \right\}$$

Just a constrained quadratic in  $\alpha_i$ :  
 $\tilde{\alpha} \in [0, C] \Rightarrow \alpha_i = \tilde{\alpha}; \tilde{\alpha} \leq 0 \Rightarrow \alpha_i = 0$   
 $\Rightarrow \frac{1}{2} q \alpha^2 - \alpha(1-p) = 0 \quad \tilde{\alpha} \geq C \Rightarrow \alpha_i = C$

$$\tilde{\alpha} = (1-p)/q, \quad q = \|x^i\|_2, \quad p = y^i \sum_{j \neq i} \alpha_j y^j \langle x^i, x^j \rangle$$

But,  $w = \sum \alpha_i y^i x^i \rightsquigarrow p = w^T x^i - \alpha_i y^i \|x^i\|_2^2$

$\Rightarrow$  We can update  $w$ , & update  $\alpha_i$  in  $O(d)$

## Gradient

Primal Gradient Descent :  $O(nd)$  time / update  
 Dual Gradient Ascent :  $O(nd)$  time / update

## Stochastic Gradient Methods

Stochastic PGD :  $O(d)$  time / update  
 Stochastic DGA :  $O(n)$  time / update

## Coordinate Methods

Stochastic Primal Coordinate Descent :  $O(d)$

Stochastic Dual Coordinate Maximization :  $O(d)$

Case 1:  $n \gg d$  : SDCM or SPGD

Case 2:  $d \gg n$  : SDGA or SPCD

★ Co-ordinate Ascent: small steps to ↑ obj a bit  
" Maximization: Completely maximize obj.

## ML primitives

- 1) What sort of action we wish to take on test data?
- 2) In what form is training data available
- 3) What sort of model we learn & how we learn it

## Primitives by type of action

- ★ Simplest primitive → Yes / No question  
Ex: Spam / Non-Spam ↗ 2 classes  
Often called **binary classification**
- ★ What if more than 2 classes?  
⇒ **Multi-class classification**  
↳ Can be seen as sending an email to one bin  
What if we want to send it to more bins?  
↳ Useful to call it a labeling / tagging problem than a classification one.  
↳ **multi-label classification**
- We assume that # & the **tags** remain consistent throughout training and testing  
↳ Can remove assumption w/ **zero-shot learning**
- ★ Assign a value / score to each data point  
↳ **Regression**  
Ex: 0 (Non-Spam) —————→ Real valued regression  
5 (Definitely Spam)

Definitely → Can change scale/range too

But, we need to modify train/test data accordingly

A few examples:

1) User vs fruit recommendation

Technique - 1: via multi-label classification

Fruits: label; likes fruit: tag them;  
User: data point

→ No. of labels might ↑↑

T-2: via regression

Scale (0-5): how much liked by user;  
(user, fruit): data point (User, ) = 4.8

2) Machine Translation

Ex: My name is Melbo →

मेरा नाम Melbo है।

T-1: Multi-Class classification

[ ] [ ] ... [ ] [ ]  
मेरा नाम Melbo है।

English sentence: data point

(My name is melbo) → data point → Classes  
Repeat till classes ↑ but this is the method used done

Primitives by form of data

\* Training data → ML algo → Model

Most popular & simplest setting of learning → supervised learning

- Variant
- ★ **Demi-supervised learning** → Some training data left unlabelled
  - ★ **Online learning** → Model given to us trained
    - Model learns ← User comes & model predicts & user added to train data afterwards
    - Promotes continuous learning
    - Train data not given to us in one go
    - No distinction b/w train/test data
  - ★ Other variants: Active Learning / Reinforcement learning

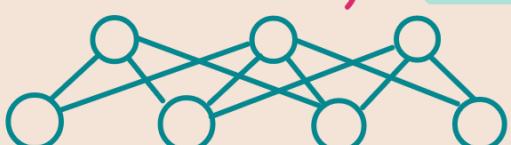
Primitives by type of model

★ **Geometric models:**

Support vector Machines (use hyperplanes),  
Ridge regression

★ **Neural models:**

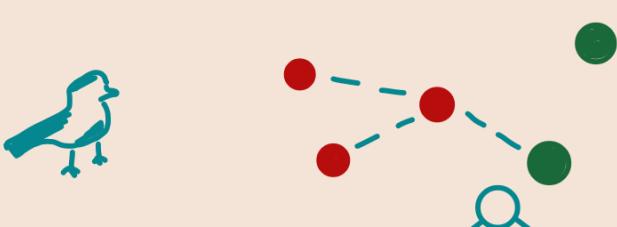
Multi-layer perceptron,  
networks, transformers



★ **Probabilistic Models:**  
Logistic regression, Bayesian models

$$P[\text{spam} \mid \text{✉}] < 0.05$$

★ **Models with memory:**



# Learning with prototypes, k-nearest neighbours, decision trees



Vectors: An ordered list of numbers

Matrices: An ordered rectangular array of numbers

## Interpreting Model Outputs (Regression)

→ Predict temperature tomorrow

↳ Want model to predict b/w [-10, 50]

↳ Ask for scores b/w [0, 1] → Scale

↳ Use sigmoid:  $\sigma(z) = \frac{1}{1+e^{-z}}$  → Shift

Example of an activation / quashing function

→ If we want sigmoid to be sharper, use  $z \rightarrow \lambda z$ ,  $\lambda > 1$

## Classification

→ Binary: Ask 1/0 {or +1/-1 in some cases} & then manually map them to labels

→ Multi-class: Ask one-hot vectors (<0, 1, 0, 0>) w/ dimensions equal to # classes  
(Can use softmax activation)

→ Multi-label: Ask multi-hot vectors

## Features of Data points

- Numerical features — As is ; age, temp, etc.
- Categorical features — 1-hot vector ; gender, blood group, etc.
- Relational features
- Many-hot vectors

## Derived features

- Bagged / Binned features

Ex: Bag of words,  $\begin{cases} \text{ESC101(A), 207(B)} \\ \text{MSD(C), 771(A), } \end{cases} \rightarrow \begin{cases} A=3 \\ B=1 \end{cases}$   
 $\begin{cases} 340(\text{A}), 220(\text{C}) \\ \end{cases} \quad C=2$

- Pooled / aggregated features

Ex: A has friends B, C, F.  $\rightarrow$  # A friends = 3  
 not friends w/ D, E, G, H  
 OR take avg & stdv to represent data

## Q Find my salary

→ Conjecture some law:

base salary  $b$ ;  $\uparrow$  by  $d$  / year of age;

$\Delta = p/q/r/s$  depending on education;  
 $M/F/I/O \rightarrow \Delta = k/l/m$

$$\rightarrow \hat{y} = b + d x_1 + p x_2 + q x_3 + r x_4 + s x_5 + k x_6 + l x_7 + m x_8 \rightsquigarrow \hat{y} = b + w^T x$$

$$x = \langle x_i \rangle_{i=1}^{i=8}, b = b, w^T = \langle d \ p \ q \ r \ s \ k \ l \ m \rangle \rightarrow \hat{y} = f(x; w, b) = b + w^T x$$

function ( $f$ ) takes data point ( $x$ ) and gives output ( $\hat{y}$ ) w/ some unknowns called the parameters (vector  $w$ , scalar  $b$ ) of the model

→ Goal of ML algos is to identify parameters

Model: form of the function  $f$ , could have been  $g = (w^T x)^b$   
 $\downarrow w^T x + b \rightsquigarrow$  Linear Model

Needs to find which params are "good":

↓ If predictions ✓, then params good  
↳ predictions match true outputs

Need to make this more mathematical

→ Common to use loss functions:  
→ let feature vector of  $i^{\text{th}}$  user:  $x^i$ ,  
true salary  $y^i$ , predicted salary  $\hat{y}^i$

→ mismatch can be measured as

$$l_i = |\hat{y}^i - y^i| \quad \text{OR} \quad (\hat{y}^i - y^i)^2$$

Absolute loss

Least squares loss

Can't take just one data point

→ Usually average of loss (across all data points in training data) is used

A duality: Data points and params operate on each other

ML model: Gives prediction  $\hat{y}^i$  on data point

$$\hat{y}^i = f(x^i; \theta)$$

$x^i$  using param values  $\theta$

Loss function: Gives loss value  $l_i$  for params

$$l_i = l(\theta; x^i, y^i)$$

Often loss value computed using predicted and true outputs  $l_i = l(\hat{y}^i, y^i) = l(f(x^i; \theta), y^i)$

Almost always, average loss value over entire training set is used

$$l(\theta) := \frac{1}{n} \sum_{i=1}^n l(\theta; x^i, y^i)$$

Use this loss as feedback to train the model



Train + test data should be similar

NEVER look at test data during training

→ Online learning obeys this because

data is first tested upon and then trained upon, NOT the other way around

- \* Models : take data points as input & use parameters to make predictions
- ML algos : use loss functions to search for good parameter values

## Geometry Basics

- Coordinates of a vector  $\equiv$  Point represented by vector
  - Stretch :  $|λ| > 1$
  - Shrink :  $0 < |λ| < 1$
  - Flip :  $λ < 0$
  - Adding two vectors
  - Subtracting " "
  - Length of a vector
- Scalar Multiplication**
- ↳ Element-wise mult.
- $u + v \rightarrow \dots \quad \text{addition}$
- $u - v \equiv u + (-1) \cdot (v)$
- Can also use || gm law
- Euclidean Length  
 $u = \langle x, y \rangle \rightsquigarrow \sqrt{x^2 + y^2}$
- Taxicab / Manhattan length  
 $u = \langle x, y \rangle \rightsquigarrow |x| + |y|$
- Norms ; Entire family of so-called  $L_p$  norms defined as :

$$\|x\|_p := (|x_1|^p + |x_2|^p + \dots + |x_d|^p)^{\frac{1}{p}}$$

\* Notice that  $L_2 \equiv$  Euclidean,  $L_1 \equiv$  Manhattan

- Lengths can be used for distances
- ↳  $d_{\text{Euclid}}(u, v) := \|u - v\|_2$ ,  $d_{\text{Man}}(u, v) := \|u - v\|_1$ .
- These notions of distances are also called metrics. Can use any of the  $L_p$  norms to define a metric.

as  $d_p(u, v) := \|u - v\|_p$

Metrics satisfy 3 nice properties:

$$d_p(u, v) = d_p(v, u)$$

$$d_p(u, u) = 0$$

$$d_p(u, v) \leq d_p(u, w) + d_p(w, v)$$

{ Symmetry }

{ Identity }

{ triangle inequality }

Measuring angles:

Dot product  $\rightarrow \langle u, v \rangle := u^T v = \sum_{i=1}^d u_i \cdot v_i$

Also,  $\langle u, v \rangle = \|u\|_2 \|v\|_2 \cos \theta \Rightarrow \theta = \cos^{-1} \left( \frac{\langle u, v \rangle}{\|u\|_2 \|v\|_2} \right)$

Obtuse angle  $\Leftrightarrow$  dot product  $< 0$

Acute angle  $\Leftrightarrow$  dot product  $> 0$

Perpendicular  $\Leftrightarrow$  dot product  $= 0$

## Application of Norms

$B_2(c, r) := \{u : \|u - c\|_2 < r\}$

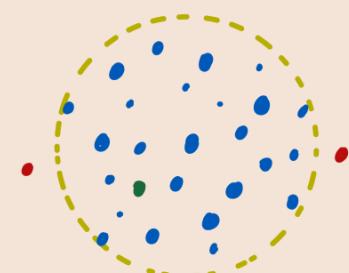
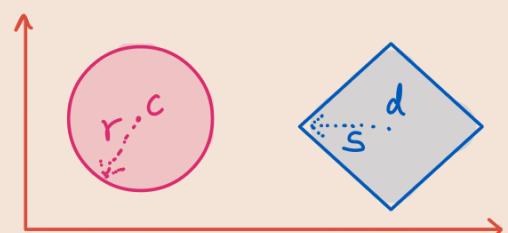
$B_1(d, s) := \{v : \|v - d\|_1 < s\}$

Anomaly detection:

Use algorithms to find smallest enclosing ball

↪ Safe if pt. inside

↪ Danger if pt. outside



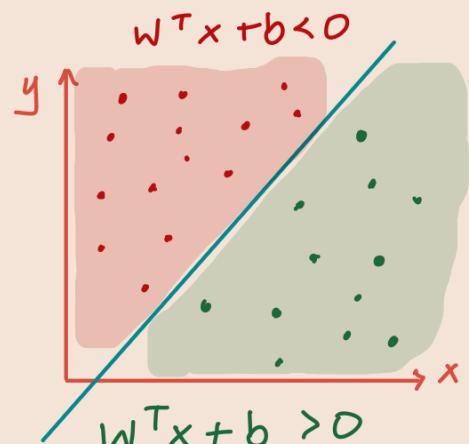
Dot product

Binary Classification

2D  $\rightarrow$  line,  $y = mx + b$

Re-write  $w^T x + b = 0$

$w = (m, -1)$ ,  $x = (x, y) \in \mathbb{R}^2$



- Both sets called **halfspace(s)**, separated by a line / **hyperplane**.
- Same trick works in higher dimensions too  $\rightsquigarrow w, x \in \mathbb{R}^d, b \in \mathbb{R}$   
The hyperplane  $\{x \mid w^T x + b = 0\}$  is called the **decision boundary**
- $\vec{w}$  is the **normal** or  $\perp$  of hyperplane  
 $\because x, y \in$  hyperplane  $\Rightarrow w^T(x-y) = 0 \Rightarrow w \perp x-y$   
 $\Rightarrow w \perp$  to all such difference vectors
- \* **bias**:  $b \uparrow \downarrow \rightarrow$  hyperplane shifts {which are  $\parallel$  to plane}  
 $b \uparrow \rightarrow$  more points classified as green  
 $b \downarrow \rightarrow$  less " " " "
- Changing  $w \rightarrow$  Plane rotates

To b or not to b  
often critical for the model to classify better

ML algos simpler if no bias term

→ Chat a bit & hide it

→ Create another dim in feature vector & fill it w/ 1 ie  $\tilde{x} = [x, 1]$

→ Learn a  $(d+1)$ -dimensional linear model w/o bias term,  $\tilde{w} = [w_1, w_2, \dots, w_d, w_{d+1}]$   
 $\tilde{w} \in \mathbb{R}^{d+1}$ . Let  $w = [w_1, w_2, \dots, w_d] \in \mathbb{R}^d$ , then  
 $\tilde{w}^T \tilde{x} = w^T x + w_{d+1}$

Thus,  $w_{d+1}$  effectively acts as a bias term

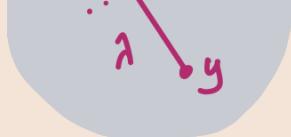
## Convex Sets



$$C \subseteq \mathbb{R}^d$$

$$\forall x, y \in C$$

$$\forall \lambda \in [0, 1] \quad \lambda x + (1-\lambda)y \in C$$



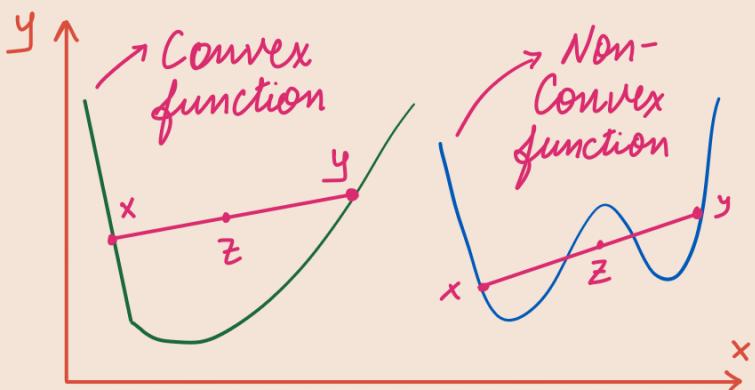
$$\begin{aligned} & \lambda \in [0, 1] \\ & z = \lambda \cdot x + (1-\lambda) \cdot y \\ \Rightarrow & z \in C \end{aligned}$$

Convex Set

Non-Convex Set

## Convex Functions

$$\begin{aligned} f: \mathbb{R}^d &\rightarrow \mathbb{R} \\ \forall x, y \in C \\ \forall \lambda \in [0, 1] \\ z &= \lambda \cdot x + (1-\lambda) \cdot y \\ f(z) &\leq \lambda \cdot f(x) + (1-\lambda) \cdot f(y) \end{aligned}$$



- ★
  - All Constant, linear functions are convex
  - Sums of convex functions are convex
  - Positive multiples " " " "
  - $g: \mathbb{R}^d \rightarrow \mathbb{R}$  convex ;  $f: \mathbb{R} \rightarrow \mathbb{R}$  convex & non-dec. then the function  $f \circ g: \mathbb{R}^d \rightarrow \mathbb{R}$  convex
  - $f: \mathbb{R} \rightarrow \mathbb{R}$  convex  $\Rightarrow h(x) = f(a^T x + b)$  convex

## Calculus Basics

Derivatives → Help us understand what happens to a function's output if its input is changed/perturbed a little

Taylor's Theorem  $\xrightarrow{\text{Corollary}}$   $f(x_0 + \Delta x) \approx f(x_0) + f'(x_0) \cdot \Delta x$



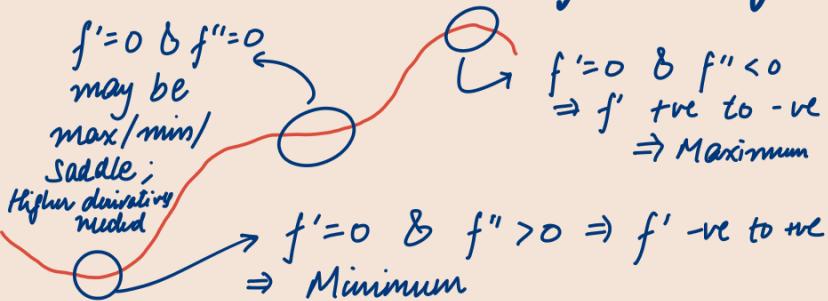
→ Holds only if  $\Delta x$  small

→  $\Delta x$  sign same as  $f'(x)$   $\Rightarrow f(x_0 + \Delta x) > f(x_0)$   
     " opposite " "  $\Rightarrow f(x_0 + \Delta x) < f(x_0)$

· Sign of  $f'(x_0)$  tells us in which dir" will function value ↑ :  $+ \equiv$  right ;  $- \equiv$  left

· Magnitude of  $f'(x_0)$  tells how much it will change

- Stationary Points**: places where derivative vanishes  
 Could be local/global minimum/maximum or a saddle point
- $f' = 0$  tells us that  $f \approx$  flat around that point  
 Very small scales
  - Can separate min. w/ max. using 2<sup>nd</sup> derivative.  
 $f''$  sign tells in which dir'  $f'$  ↑  
 Magnitude " how much  $f'$  changes"
  - No general way of telling if a max/min is local or global



- Some rules:  $(f+g)' = f' + g'$ ;  $(af)' = a \cdot f'$ ,  $a \in \mathbb{R}$  const.  
 $(f(g))' = f'(g) \cdot g'$ ;  $(fg)' = fg' + f'g$ ;  $(f/g)' = (f' \cdot g - f \cdot g')/g^2$   
 $\{ f(x), g(x) : \mathbb{R} \rightarrow \mathbb{R} \}$

## Multivariate Derivatives

a.k.a. Gradients

Consider  $f(x, y) : \mathbb{R}^2 \rightarrow \mathbb{R}$

Trick: Convert the problem into analyzing  $\mathbb{R} \rightarrow \mathbb{R}$  functions

→ Fix all variables except one then diff. w.r.t. that

Ex:  $\frac{\partial f}{\partial x}$  at  $(x_0, y_0)$ : derivative of  $f$  w.r.t.  $x$   
 keeping  $y_0$  const

↳ Sign tells us  $\uparrow\downarrow$  if  $\Delta x > 0$  ( $y = y_0$  fixed)

↳ Mag. tells how sharply  $f$  changes on  $\Delta x$

Similarly  $\frac{\partial f}{\partial y}$

Vector  $\left( \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right) := \nabla f(x, y)$  ← Gradient of  $f$

Co-ordinate-wise derivatives arranged in the form of vector

$f: \mathbb{R}^d \rightarrow \mathbb{R}$  then  $\nabla f(x) = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \frac{\partial f}{\partial x_3}, \dots, \frac{\partial f}{\partial x_d} \right)$

Stationary points  $\rightarrow \nabla f = 0$

## Steepest ascent

$$f(x^0 + t) \approx f(x^0) + t^T \nabla f(x^0)$$

$f: \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $t = (t_1, t_2, \dots, t_d) \in \mathbb{R}^d$ ,  $t$  is small:

$$f(x^0 + t) \approx f(x^0) + \sum_{i=1}^d t_i \cdot \frac{\partial f}{\partial x_i}$$

Claim: The dir<sup>n</sup> of the gradient vector offers the biggest  $\uparrow$  in function value out of all directions

## Steepest descent

For proof, use  $f(x^0 + t) \approx f(x^0) + t^T \nabla f$   
 $t^T \cdot \nabla f = \|t\|_2 \|\nabla f\|_2 \cos \theta = \langle t, \nabla f \rangle$   
& then fix  $\|t\|_2 = E$

Claim: The dir<sup>n</sup> opposite to gradient vector offers the biggest  $\downarrow$  in function value out of all directions

Some rules:  $f, g: \mathbb{R}^d \rightarrow \mathbb{R}$  &  $s: \mathbb{R} \rightarrow \mathbb{R}$

Sum Rule:  $\nabla(f(x) + g(x)) = \nabla f(x) + \nabla g(x)$

Scaling "":  $\nabla(a \cdot f(x)) = a \cdot \nabla f(x)$ ,  $a \in \mathbb{R}$  constant w.r.t.  $x$

Product "":  $\nabla(f(x) \cdot g(x)) = g(x) \cdot \nabla f(x) + f(x) \cdot \nabla g(x)$

Quotient "":  $\nabla(f(x)/g(x)) = [g(x) \cdot \nabla f(x) - f(x) \cdot \nabla g(x)]/[g(x)]^2$

Chain "":  $\nabla(s(f(x))) := (s \circ f)'(x) = s'(f(x)) \cdot \nabla f(x)$

## Some identities:

$c \in \mathbb{R}$  constant that doesn't vary w/  $x$   $\nabla(c) = 0$

$a \in \mathbb{R}^d$  "vector" "  $\nabla(a^T x) = a$

$$\begin{array}{ll} \cdot A \in \mathbb{R}^{d \times d} \text{ "matrix"} & \nabla(x^T A x) = (A + A^T)x \\ \cdot a \in \mathbb{R}^d \text{ "vector"} & \nabla(\|x - a\|_2) = \frac{x - a}{\|x - a\|_2} \end{array}$$

$$\nabla(x^T A x) = (A + A^T)x$$

(= 2Ax if A symmetric)

Ex:  $w^T x + b$  on  $\{(x^i, y^i)\}_{i \in [n]}$ ;  $b, y^i \in \mathbb{R}$ ;  $w, x \in \mathbb{R}^d$   
w/  $\ell(w, b) = \frac{1}{n} \sum \|y^i - (w^T x^i + b)\|^2 = \frac{1}{n} \sum ((w^T x^i + b) - y^i)^2$

- Simplify by hiding bias  $b$ :  $\tilde{x}^i = [x^i \ 1]$ ,  $\tilde{w} = [w, 1]$
- $\ell(\tilde{w}) := (\tilde{w}^T \tilde{x}^i - y^i)^2 = \underbrace{(\tilde{w}^T \tilde{x}^i)^2}_{\nabla = \nabla(\tilde{w}^T \tilde{x}^i \tilde{x}^{i\top} \tilde{w})} + \underbrace{(y^i)^2}_{\nabla = 0} - \underbrace{2y^i \cdot \tilde{w}^T \tilde{x}^i}_{\nabla = 2y^i \tilde{x}^i}$
- $\nabla \ell(\tilde{w}) = 2\tilde{x}^i \tilde{x}^{i\top} \tilde{w}$

★ Let  $v = [v_1 \ v_2 \ \dots \ v_d] \in \mathbb{R}^d$  then

- $v \cdot v^T = d \times d$  square symmetric matrix  
whose  $(i, j)^{\text{th}}$  entry is  $v_i \cdot v_j$
- $v^T \cdot v = \|v\|_2^2 \in \mathbb{R}^+$  { Squared Euclidean length }

- $\nabla \ell(\tilde{w}) = 2/n (\sum \tilde{x}^i \tilde{x}^{i\top}) \tilde{w} - 2/n (\sum y^i \cdot \tilde{x}^i)$
- As  $\ell(\tilde{w})$  diff., minima must be a stationary point  
 $\Rightarrow \nabla \ell(\tilde{w}) = 0 \Rightarrow (\tilde{X}^T \tilde{X}) w = \tilde{X}^T y$

where  $\tilde{X} = \begin{bmatrix} \tilde{x}^1 \\ \tilde{x}^2 \\ \vdots \\ \tilde{x}^n \end{bmatrix} = \begin{bmatrix} x^1 & 1 \\ x^2 & 1 \\ \vdots & \vdots \\ x^n & 1 \end{bmatrix} \in \mathbb{R}^{n \times (d+1)}$  &  $y = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^n \end{bmatrix} \in \mathbb{R}^n$

- If  $(\tilde{X}^T \tilde{X}) \in \mathbb{R}^{(d+1) \times (d+1)}$  invertible, then  
the solution must be  
If not invertible,  
then techniques  
such as regularization are used

$$\tilde{w} = (\tilde{X}^T \tilde{X})^{-1} (\tilde{X}^T y)$$

---

Subgradients and  
Subdifferentials

→ for convex diff. funcs, gradient defines a tangent hyperplane at every point & the function must lie above this plane

→ Subgradients exploit and generalize this property to even points of non-differentiability.

\* Convex functions lie above all tangents  
 $(\forall x_0) \quad f(x) \geq \nabla f(x_0)(x - x_0) + f(x_0) \quad \forall x$

Turn definition around and say that gradient is a vector  $g$  so that the hyperplane  $0 = g^T(x - x_0) + f(x_0)$  is tangent to  $f$  at  $x_0$ .

**Subgradients:** Given a (possibly non-diff. but convex) function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  & point  $x^0$ , any vector  $g$  that satisfies

$$f(x) \geq g^T(x - x^0) + f(x^0)$$

is called a subgradient of  $f$  at  $x^0$

**Subdifferential:** The set of all subgradients of  $f$  at a point  $x^0$  is known as the subdifferential of  $f$  at  $x^0$  & denoted by  $\partial f(x^0)$

• Function may have multiple subgradients at  $x^0$  but if it is diff. at  $x^0$  then the gradient is the only subgradient at  $x^0$ .

$$\partial f(x^0) := \{ g \mid f(x) \geq g^T(x - x^0) + f(x^0) \}$$

Subgradient Calculus



- Scaling Rule:  

$$\partial(c \cdot f(x)) = c \cdot \partial f(x)$$

$$= \{c \cdot v : v \in \partial f(x)\}$$
  - 
  - Sum Rule:  

$$\partial(f+g)(x) \subseteq \partial f(x) + \partial g(x) = \{u+v : u \in \partial f(x), v \in \partial g(x)\}$$
  - Chain Rule:  

$$\partial f(a^T x + b) = \{c \cdot a : c \in \partial f(a^T x + b)\}$$

$$\quad \quad \quad \{ \nabla f(a^T x + b) = f'(a^T x + b) \cdot a \}$$
  - Max Rule:  $h(x) = \max \{f(x), g(x)\}$ 
    - $f(x^o) > g(x^o) \Rightarrow \partial h(x^o) = \partial f(x^o)$
    - $f(x^o) = g(x^o) \Rightarrow h(x^o) = \begin{cases} \lambda u + (1-\lambda)v : \\ u \in \partial f(x^o), v \in \partial g(x^o), \lambda \in [0,1] \end{cases}$

$$\text{Ex: } \text{hinge}(x) = \max(1-x, 0)$$

$$\downarrow = \max(f(x), g(x))$$

diff.  $\forall x \in \mathbb{R} \setminus \{1\} \Rightarrow \partial \ln|x| = \ln|x| \text{ if } x \neq 1$

$$\text{At } x=1 \rightarrow f(x) = 1-x \rightarrow \partial^0 f(x) = f'(x) = -1$$

$$g(x) = 0 \Rightarrow \partial g(x) = 0$$

$$\therefore \partial \lambda_{\text{mige}}(1) = \{ \lambda \cdot (-1) + (1-\lambda)(0) = -\lambda : \lambda \in [0, 1] \} = [-1, 0]$$

Applying chain rule, we get,

$$\text{change}(y^i, \langle w, x^i \rangle) = [i - y^i \langle w, x^i \rangle]_+$$

Need  $v^i \in \partial L_{\text{ hinge}}(y^i, \langle w, x^i \rangle)$

$$v^i = \begin{cases} 0, & y^i \langle w, x^i \rangle > 1; \\ c \cdot y^i \cdot x^i, & c \in [-1, 0], \quad y^i \langle w, x^i \rangle = 1 \\ -y^i \cdot x^i, & y^i \langle w, x^i \rangle < 1; \end{cases}$$

## BTS in GD for SVM

$$f(w) = \frac{1}{2} \|w\|_2^2 + C \cdot \sum_{i \in [n]} [1 - y^i \cdot w^T \cdot x^i]_+$$

{ ignore bias  $b$  }

$$\nabla f(w) = w + C \cdot \sum_i g^i y^i \cdot x^i, \quad g^i \in \partial \ell_{\text{ hinge}}(y^i \cdot w^T x^i)$$

$$w^{new} = w - \eta \cdot \nabla f(w) = (1-\eta) \cdot w - \eta C \cdot \sum_i g^i y^i \cdot x^i$$

8  $w^{new} = w - \eta \cdot \nabla f(w) = (1-\eta) \cdot w - \eta C \cdot \sum g_i y_i x_i$

Assume  $n=1$  for the sake of understanding.

$$w^{new} = (1-\eta) \cdot w - \eta C \cdot g' \cdot y' \cdot x'$$

Small  $\eta$ :  $(1-\eta)$  large  $\Rightarrow$  don't change  $w$  too much

Large  $\eta$ : Change  $w$  as much as gradient dictates

$\hookrightarrow w$  did well on  $(x', y')$ , say  $y' \cdot w^T x' > 1 \Rightarrow g' = 0$

$\rightarrow$  did bad, say  $y' \cdot w^T x' < 0$ , then  $g' = -1$

$$\Rightarrow w^{new} = (1-\eta) \cdot w + \eta C \cdot y' \cdot x'$$

$$\Rightarrow y' \cdot (w^{new})^T \cdot x' = (1-\eta) y' \cdot w^T x' + \eta C \cdot \|x\|_2^2$$

$w$  too much  
don't change

$w^{new}$  gets  
better margin than  $w$

## Stochastic Gradient Method

$$\nabla f(w) = w + C \cdot \sum g^i y^i \cdot x^i, \quad g^i \in \nabla \text{loss}_i(y^i \cdot w^T x^i)$$

$O(d)$   
time  
only

Take a random  
data point instead  
of update w.r.t. it

Calculating each  $g^i$  takes  
 $O(d)$  time  $\Rightarrow n$  data points  
Total  $O(nd)$  time

$$\nabla f(w) = w + C \cdot g^{i_t} y^{i_t} x^{i_t}$$

Might have to do more  
steps than in GD but each  
step cheaper, hence faster

## Mini-Batch SGD

If data diverse, the stochastic gradient may vary quite a lot depending on which random data point chosen.  $\rightsquigarrow$  Called variance

$\Rightarrow$  Choose a bunch ( $B$  to be precise) of data points and update w.r.t. them

$\hookrightarrow O(B \cdot d)$  time per update

$$\nabla f(w) \approx w + C \cdot \sum_{b=1}^B g^{i_b} y^{i_b} x^{i_b}$$

- \* If  $B = n$ , then MBSGD = GD, and if  $B = 1$ , then MBSGD = SGD

## Coordinate Descent

Similar to GD except only one coordinate is changed in a single step

$\hookrightarrow \min_{x \in \mathcal{C}} f(x) \text{ s.t. } x \in \mathcal{C} \text{ w/ } \nabla_j f(x) = \frac{\partial f}{\partial x_j} \rightsquigarrow j^{\text{th}}$  partial derivative

$\nabla f(x)$

- CCD: Choose coordinate cyclically,  $j_t = 1, 2 \dots d, 1, 2 \dots d, \dots$
- SCD: Choose  $j_t$  randomly
- Block CD: Choose a small set of coordinates to update
- Randperm: Permute randomly and choose.  
Once the list is over, choose a new randperm.
- Sometimes, we are able to optimize completely along a given variable, called coordinate minimization (CM)

For  $t = 0 \dots$

- Select  $j_t \in [d]$
- $x_{j_t}^{t+1} \leftarrow x_{j_t}^t - \eta_t \cdot \nabla_{j_t} f(x^t)$
- $x_j^{t+1} \leftarrow x_j^t \quad \forall j \neq j_t$
- Repeat until convergence

## Issues with GD

- Initialization : Bad initialization  $\rightarrow$  slow convergence unless step lengths nice [convex func]
- [Non-convex func]: Bad init  $\rightarrow$  Stuck at saddle points
- Random restarts most common solution to overcome
- For some nice non-convex problems, we know very good ways to provably init close to global optimum
- Deciding Convergence:
  - Small dist. of global / local optima
  - Not making much progress,  $\|w^{t+1} - w^t\| \rightarrow 0$
  - GD stops on reaching a stationary point  
 ↳ Progress stopped w/o a global optimum
  - Few heuristics to stop GD:
    - If gradient vectors too small:
    - If obj. function already acceptably small:
- Detecting Convergence:

### M1: Tolerance Technique

- Pre-decided tolerance  $\epsilon$ , if  $f(w^t) < \epsilon$ , stop

### M2: 0<sup>th</sup> order technique

- If func. value doesn't change much, stop (or tune learning rate)

$$|f(w^{t+1}) - f(w^t)| < \tau \quad \text{OR} \quad |w^{t+1} - w^t| < \zeta$$

M3: 1<sup>st</sup> order technique

- If gradient too small,  $\|\nabla f(w^t)\|_2 < \delta$

M4: Cross validation technique

- Test current model on validation data

Other techniques like primal-dual are usually infeasible for large scale ML problems

→ Choosing Step length

- For nicely behaved convex functions, have formulae for step length
- Set  $\eta_t = C/\sqrt{t}$  or  $\eta_t = C/t$ ,  $C$  is a hyperparameter. Idea is to choose  $\eta_t \rightarrow 0$  (diminishing) and  $\sum \eta_t \rightarrow \infty$  ( $\infty$  travel)
- For "nice" convex functions,  $\in$  convergence in just  $O(\log(\gamma\epsilon))$  steps.

\* Powerful but expensive technique:

Newton method:

$$g^t = (\nabla^2 F(w^t))^{-1} \nabla F(w^t)$$

Autotunes the step length so we may use

$$w^{t+1} \leftarrow w^t - g^t$$

- For nice convex func,  $\in$  convergence in just  $O(\log \log(\gamma\epsilon))$  steps
- But, computation of Hessian  $\nabla^2 F(w^t)$  is expensive but can be approximated using a diag or a low-rank matrix
- For not so well-behaved convex & non-convex functions,  $\exists$  several heuristics but no guarantee for them to work
  - Line-search: Find best  $\eta_t$  every time
  - $\eta_t = \arg \min_{\eta \geq 0} f(w^t - \eta \cdot g^t)$
  - Armijo rule:

- Start w/  $\eta_t$ , if it doesn't reduce obj. func. sufficiently,  $\downarrow \eta_t$  & repeat.
- Momentum methods like NAG, Adam, which infuses prev. gradients into cur. grad.
- Use diff. step length for each dimension of  $w$  (e.g. Adagrad) where  $\eta_t$  replaced by a diag. matrix  $E^t$ ,  $w^{t+1} \leftarrow w^t - E^t \cdot g^t$

## Loss Functions

- Used in ML to ask the model to behave a certain way.
- Penalize undesirable behaviour** and optimizers like SGD/SDCA then give us a model w/ desirable behaviour
- Ex. hinge loss function penalizes if it either misclassifies a data point or else classifies it correctly w/ an insufficient margin.
- Might take avg. or sum of all data points.
- Ex:

\* Hinge

\* Squared Hinge

\* Logistic

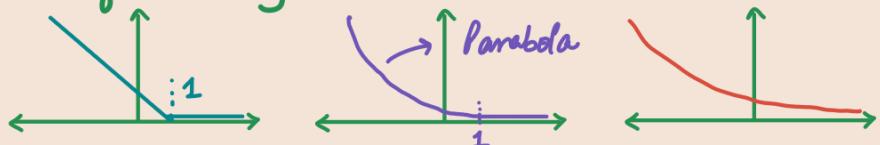
Differentiable

Related to cross-entropy

$$\text{hinge}(y_i \cdot w^\top x_i) = [1 - y_i \cdot w^\top x_i]_+$$

$$\text{sq-hinge}(y_i \cdot w^\top x_i) = [1 - y_i \cdot w^\top x_i]_+^2$$

$$\text{logistic}(y_i \cdot w^\top x_i) = \ln(1 + e^{-y_i \cdot w^\top x_i})$$



**Logistic Regression:** Obtaining a model using logistic loss function, i.e.,  $\min_{w, b} \frac{1}{2} \|w\|_2^2 + C \sum_{i=1}^m \ln(1 + e^{-y_i w^\top x_i})$

- Can be solved in primal using GD/SGD/MB
- Can be solved in dual using SDCA
  - ↳ Dual derivation not that straightforward
- Doesn't perform regression

## Regression Problems

- Predict a real value
- Training data looks like  $\{(x^i, y^i)\}_{i=1}^n$   
where  $x^i \in \mathbb{R}^d$ ,  $y \in \mathbb{R}$

## Loss functions for Regression Problems

- Can use linear models,  $(w, b)$ ,  
and predict score for test  
data point  $x^t$  as  $\hat{y}^t = w^T x^t + b$

- Absolute Loss:  $l_{abs}(\hat{y}, y) = |\hat{y} - y|$  Bad if  $\hat{y} \gg y$   
bad if  $\hat{y} \ll y$
- Squared Loss:  $l_{sq}(\hat{y}, y) = (\hat{y} - y)^2$  differentiable
- Vapnik's  $\epsilon$ -insensitive Loss:

*if slightly bad, don't punish, else penalized as a sq. func. Also, diff.*

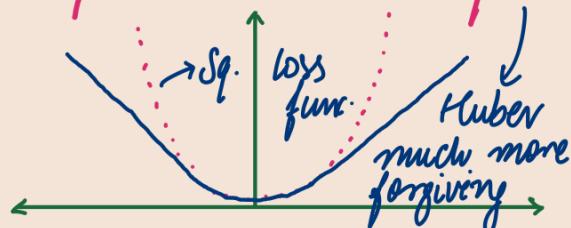
$$l_\epsilon(y, \hat{y}) = \begin{cases} (y - \hat{y} - \epsilon)^2, & \hat{y} < y - \epsilon \\ 0, & \hat{y} - y \in [-\epsilon, \epsilon] \\ (y - \hat{y} + \epsilon)^2, & \hat{y} > y + \epsilon \end{cases}$$

*slightly bad: sq. loss, else abs. loss*

- Huber Loss:  $l_h(y, \hat{y}) = \begin{cases} (\hat{y} - y)^2, & |\hat{y} - y| \leq \delta \\ \delta \cdot |\hat{y} - y|, & |\hat{y} - y| \geq \delta \end{cases}$

- \* Loss functions have to be chosen acc to needs of the problem.

- \* Huber Loss → Some data points are corrupted
- \* Hinge / Cross-entropy for binary classification
- \* Squared for regression



## Ridge Regression

{Ignore  $b$ }

$$\underset{w \in \mathbb{R}^d}{\text{arg min}} \frac{\lambda}{2} \|w\|_2^2 + \sum_{i=1}^n (w^T x^i - y^i)^2$$

Can rewrite with  $X \in \mathbb{R}^{n \times d}$ ,  $y \in \mathbb{R}^n$

First order optimality

Gradient:  $\lambda \cdot w + 2x^\top (Xw - y)$

Must vanish  $\rightarrow w = (2x^\top X + \lambda \cdot I)^{-1}(2x^\top y)$

Takes  $O(d^3)$  time

to invert the matrix  $\rightarrow$  May use (S)GD

\* Customary in regression to write optimization problem a bit differently than before

$$\underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \frac{\lambda}{2} \|w\|_2^2 + C \cdot \sum (w^\top x^i - y^i)^2$$

↳ Can use dual methods here too but not straightforward

\* Can't apply 1<sup>st</sup> Order Optimality w/ fancier loss functions (Vapnik's), need to use SGD, SDCA, etc.

\* Ridge Regression uses least squares loss & L<sub>2</sub> regularization

BTS: GD for Ridge

Regression

$$f(w) = \frac{\lambda}{2} \|w\|_2^2 + \sum_{i=1}^n (w^\top x^i - y^i)^2$$

$$\nabla f(w) = \lambda w + 2 \sum (w^\top x^i - y^i) \cdot x^i$$

$$w^{\text{new}} = w - \eta \cdot \nabla f(w) = (1 - \eta \lambda) \cdot w - 2\eta \sum (w^\top x^i - y^i) \cdot x^i$$

Assume  $n=1$

↳  $\eta \downarrow \Rightarrow (1 - \eta \lambda) \uparrow \Rightarrow$  Don't change  $w$  too much

•  $w$  does well  $\Rightarrow y' = w^\top x' \Rightarrow w^{\text{new}} = (1 - \eta \lambda) w$

•  $w$  does badly {say  $w^\top x' \gg y$ } then

{similarly for  $w^\top x' \ll y$ }  $w^{\text{new}} = (1 - \eta \lambda) w - 2\eta g \cdot x'$ ,  $g \gg 0$

$$(w^{\text{new}})^\top \cdot x' = (1 - \eta \lambda) w^\top \cdot x' - 2\eta g \|x'\|_2^2 < w^\top \cdot x'$$

\* GD merely tries to get models that behave "good"

Regularization

• An umbrella term used in ML to describe a family of terms taken to prevent

- whole family of steps taken to prevent ML algos. from suffering from problems in data
- Help offer stable behaviour
  - If data perfectly clean, no need for regularization
  - Usually involves its own hyperparameters that need to be tuned using data itself { using validation techniques }
  - Regularization can be "somewhat" data independent
  - Sometimes, makes the problem well-posed
  - How to do regularization often decided w/o looking at data

Ex.  $f(w) = \frac{1}{2} \|w\|_2^2 + C \cdot \sum [1 - y^i w^T x^i]_+$

$$f(w) = \frac{1}{2} \|w\|_2^2 + \sum [1 - y^i w^T x^i]^2$$

(squared) L2 regularizer

In regression, ensures uniqueness of solution  $\rightarrow X^T X$  non-invertible

In binary classification, encourages large margin

$\infty$  soln  $\Leftrightarrow \lambda = 0 \Leftrightarrow \lambda > 0$  ensures unique solution

- \*  $\lambda$  regularizer essentially tells the optimizer to not just blindly return a model that does well on data, but rather one that does well and is simple. [ L2 regularizer  $\rightarrow$  Simple if small L2 norm ]
- In bin. class.  $\rightarrow$  Simple models also had large margins; Don't over-regularize
- Large value of regularization  $\rightsquigarrow$  Very useless model
- Key is moderation  $\rightarrow$  Use validation
- Other considerations: How noisy is the data?  
Can afford to  $\Leftarrow$  Data  $\uparrow$   $\xleftarrow[\text{Rule of Thumb}]{}$  How much data?  
+ regularization
- Other regularizers:

LASSO:  $f(w) = \lambda \|w\|_1 + \sum (w^T x^i - y^i)^2$   
 L1-reg SVM:  $f(w) = \|w\|_1 + C \cdot \sum [1 - y^i w^T x^i]_+$

→ Prefers sparse vectors  
 Faster at test time → Often make coordinates  $\approx 0$ ,  $= 0$  to save space  
 → Consume less memory

- Popular in high dimensional problems,  $d \approx 1$  million
- $\because$  L1 non-diff., need to use subgradient methods

Although, proximal gradient descent does much better

### ★ Regularization by Early Stopping

- Before it has solved the optimization problem fully
- Automatically prevents the model from fitting the data too closely (good if noise/outliers exist)
- Often happens implicitly w/ complex optimization problems
- May get an overregularized model that does not fit data at all if stopped too early

### ★ Reg. by adding noise

- Add controlled noise so model learns to perform well despite noise
- Most well-known practice — Dropout in DL
- \* Or use subset of data that seems clean
- \* Or use corruption-aware loss functions (e.g. Huber)
- \* Or use only informative, non-noisy features out of  $d_{\text{total}}$
- Called sparse recovery
- \* Or by imposing convex constraints
- \* Or by imposing non-convex constraints

$$\underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \|Xw - y\|_2^2 \quad \text{s.t. } \|w\|_1 \leq R$$

Convex set {L<sub>2</sub> ball}

$$\underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \|Xw - y\|_2^2 \quad \text{s.t. } \|w\|_0 \leq s$$

Sparse recovery

$\|v\|_0 = \# \text{ Non-zero entries}$

Forces  $s \ll d$  features

# Multiclass Classification aka multiclassification

- Can be solved using linear models → reduce to binary classification
  - 
  - If there are  $C$  classes, train  $C$  linear models, each trained to identify one class
  - Called the OVA (one v/s all) method
  - At test time, ask all  $C$  models



OVA

- \* Create  $C$  binary classification datasets
    - For each  $c \in [C]$ , create the dataset:  $y^{i,(c)} = \begin{cases} 1 & : y^i = c \\ -1 & : y^i \neq c \end{cases}$
    - Learn a model to distinguish data points in class  $c$  from those not in class  $c$ :  $\hat{w}^c = \arg \min_w \sum_{i=1}^n l(y^{i,(c)}, \langle w, x^i \rangle)$
    - At test time, predict the class whose model gives the test point the highest score:  $\hat{y}^t = \arg \max_{c \in [C]} \langle \hat{w}^c, x^t \rangle$
    - Can introduce margin here as well
    - Demand that if true class of data point  $x$  is  $c^*$ ,  $\langle w^{c^*}, x \rangle \geq \langle w^c, x \rangle + 1 \quad \forall c \neq c^*$
    - Introducing slack forms an optimization problem
    - Can rewrite this in terms of the

## Crammer-Singer Loss

where  $\ell_{CS}(y, \{\eta_c\}_{c=1}^C) = [1 + \max_{c \neq y} \eta_c - \eta_y]_+$

Hinge Logistic  $\xrightarrow{\text{Multiclassification}}$  Crammer-Singer Loss  
 Softmax

$$\eta_c := \langle w^c, x \rangle, \quad \ell_{SM}(y, \{\eta_c\}_{c=1}^C) = -\ln\left(\frac{e^{\eta_y}}{\sum_{c=1}^C e^{\eta_c}}\right)$$

Softmax encourages  $\eta_y$  to be the largest of all  $\eta_c$   
 $\ell_{SM} \rightarrow 0$  if  $\eta_y > \text{all } \eta_c$   
 ↳ Diff unlike CS loss  $\because \sum e^{\eta_c} \approx e^{\eta_y}$  &  $\ln x \rightarrow 0$  as  $x \rightarrow 1$

## \* Other techniques $\rightarrow$ Output Codes, Decision Trees

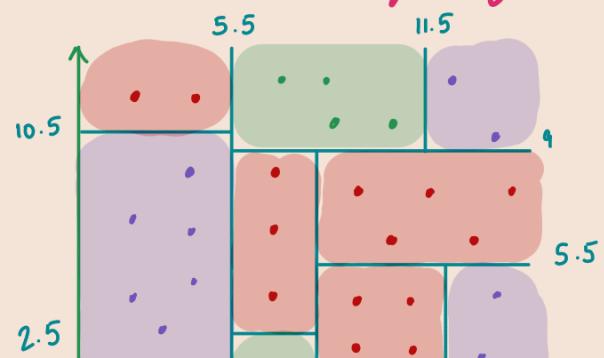
- Convert multiclassifications into regression problems very popular, especially if  $C \gg 1$
- Represent  $c \in [C]$  as a  $k$ -dim vector  $v_c \in \mathbb{R}^k$
- Solve  $k$  regression problems  $\begin{cases} \text{Find } v_1 \text{ for } c_1 \\ \text{Find } v_2 \text{ for } c_2 \dots \end{cases}$
- At test time, predict  $k$  numbers, see which vector out of  $v_1, v_2, \dots$  is it closest to
- $k$  small for speedup but can't have very small  $k$ .  $k > 1$  — Account for regression mistakes

## Decision Trees

- \* Can handle non-numeric features easily
- \* Popular in ML - classification, recsys
- \* Extremely fast at predicting \* Easy to interpret
- \* Model size can be large
- \* Can give good train w/ bad test (overfitting)

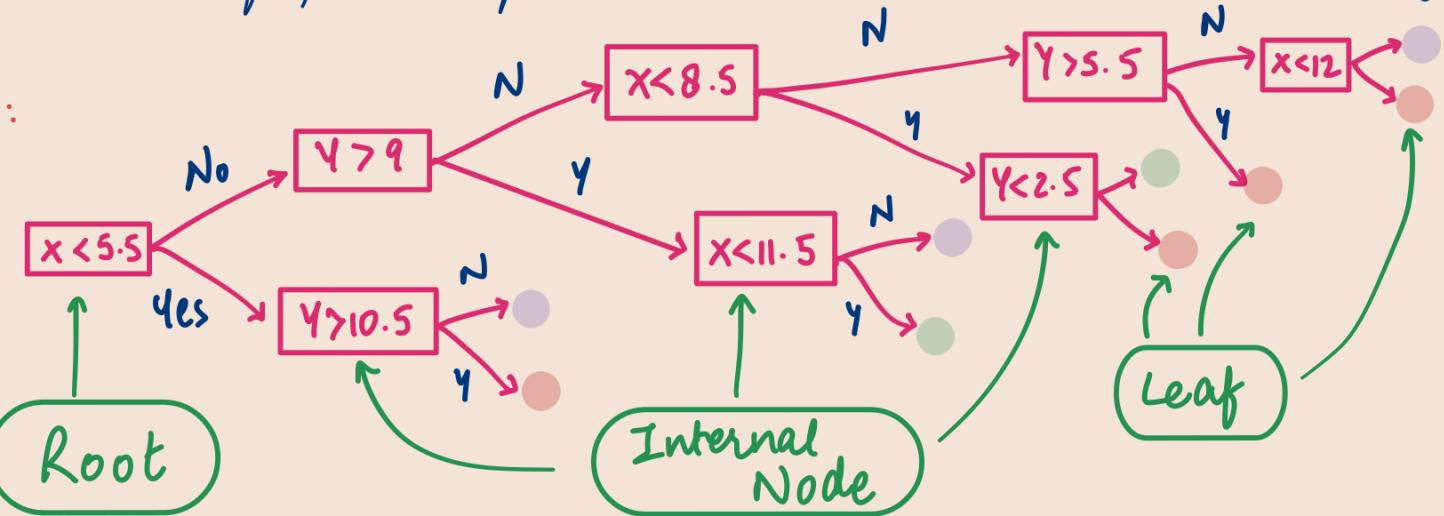
Algo: Repeatedly partition feature space  $\mathbb{R}^d$

For test data point, much faster to find out which



- partition "is it in"
  - Can consider assigning the label of other points in partition
  - Sometimes, nearest neighbours lie in other partitions
  - Rare if partitions fine enough
- DT may be wrong about points at the boundary

Ex:



Balanced DT: All leaf nodes at same depth from root  
 Imbalanced  $\rightarrow$  Bad

$\hookrightarrow$  Very poor prediction accuracy

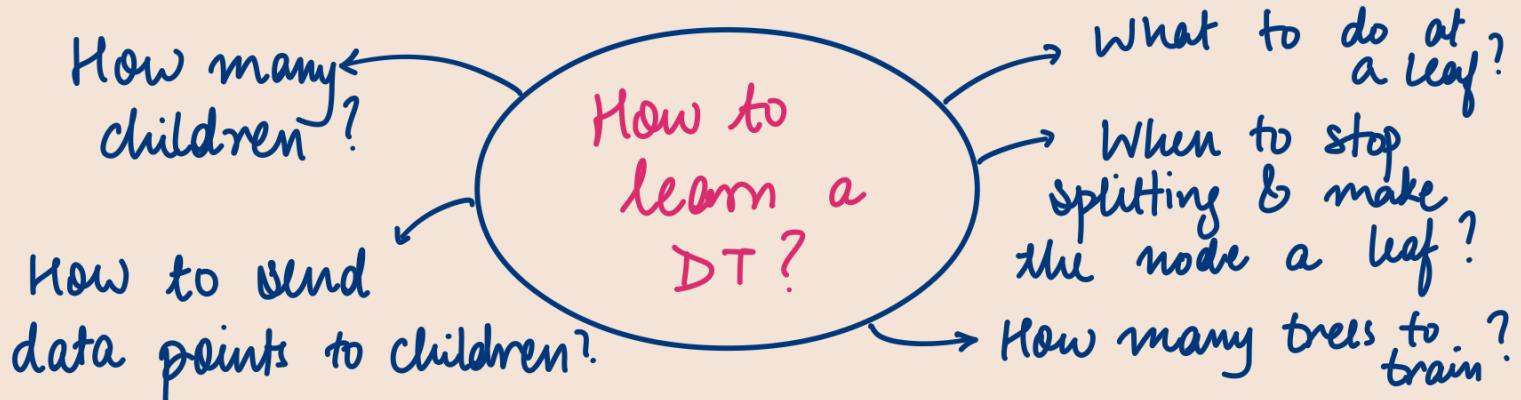
Can take as long as kNN to predict

Imagine DT: chain of nodes, n data points  
 $\Rightarrow O(n)$  chain  $\Rightarrow$  some predictions take  $O(n)$  time

Balanced DT  $\rightarrow O(\log n)$  time

can ↑ # of children/node  $\hookrightarrow$  May prune the tree to make it shallower

Regression w/ DT  $\rightarrow$  Can use avg. score of a leaf node to predict scores for test data



What to do at leaf?



11

What is it in leaf:

- can take **any** (complicated)

Another ML algo?

**action** at a leaf

For speed, keep it simple → Constant prediction

cheapest thing:  
Store majority color  
at a leaf.  
More informative (and  
expensive) thing:  
Store how many  
training points of each  
color reached that leaf.

- Such a DT will encode a piecewise const. prediction function

## ★ How to split a node?

- Splitting a node is a **classification** problem in itself. Binary if 2 children, multiclass otherwise

In principle, can use any ML classification technique (even a deepnet). But in practice, simple ML algo like linear used to retain speed → Send data point rapidly to a leaf.

One goal of DT is to speed up kNN prediction time. Thus, node splitting algo should be **superfast**, otherwise, may as well use NN.

- Often just a **single** feature is chosen and node is split based on that.
- Such simple classifiers often called decision stumps

How to find this **feature**? Choose the one that makes the **child nodes** created as **pure** as possible

**Child node pure**: if it contains

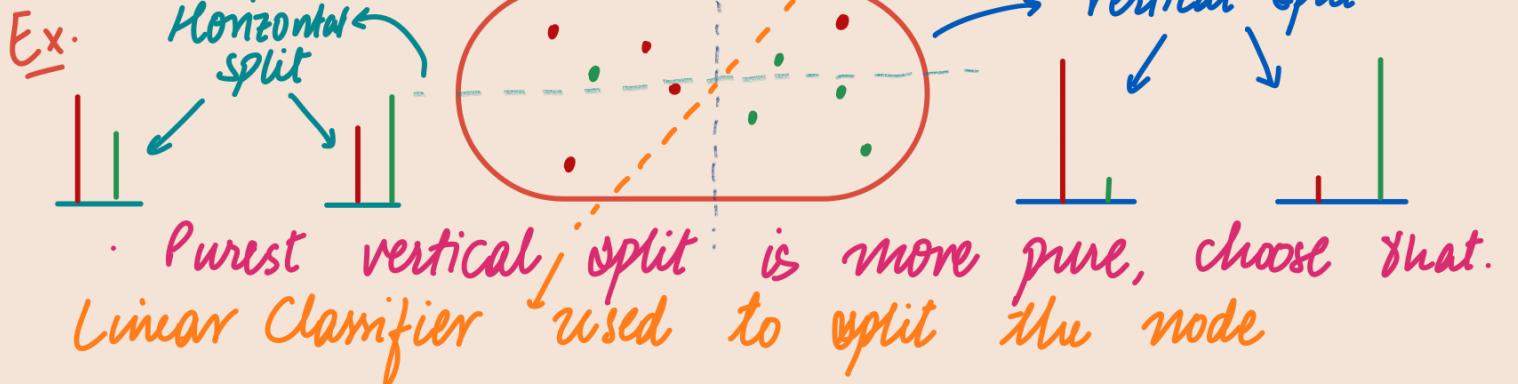
training data of only one class

Convenient because can be made leaf directly

## ★ Notion of purity?

- Various: **entropy** / Gini index for classification, **variance** for regression problems
- Making sure that split is balanced is also important, but this is often tricky.

Vertical split



## \* Pruning Strategies

- Stop if node is almost pure
- Stop if all features exhausted, avoid using a feature twice on a path, limits depth to atmost  $d$  (# of dimensions)
- Can stop if node is ill-populated (few training points)
- Can also (over)grow a tree and then merge nodes to shrink it
- Merge two leaves and see if it worsens performance on validation set — rinse & repeat
- Always use a validation set to make these decisions, never the test set

## Probability Basics

- How frequently does something happen
- Way of measuring the amount of uncertainty in a statement, Way of assigning positive scores in a way so that two scores can be easily compared

Sample Space : Denotes an exhaustive enumeration of all possible outcomes that have either happened or could happen (even if extremely unlikely)

→ Often makes the space oo in size in real settings

Event : anything a description of which tells about

Event: Simply a description of useful facts about an outcome.

- ML can be used to tell us how frequently an event occurs / if one event is more likely than other
- Tell us how confident is the ML algorithm while giving the above replies
- An event may choose to be precise about certain aspects and neglect others.

Random Variable: A way to express useful facts about events as numbers so that we can do math.

- RVs can be categorical or numerical:
  - Categorical:  $X=1$  if female, 2 male, 3 trans
  - Numerical (Discrete):  $Y$  = age of person in years
  - Numerical (Continuous):  $Z$  = no. of sec. spent on website
  - Indicator:  $W=1$  if purchase, 0 if no purchase
- Can arrange many RVs as vectors too
- Features  $\leftrightarrow$  R.V.

## Probability Distribution

- Given an event, it can tell us how likely it is  $\rightsquigarrow$  can also compare two events
- Generate a sample outcome

## Example

- Sample Space: {R, G, B}  $\times$  [6]
- $P[R] = 14/24$ ,  $P[B] = 4/24$ ,  $P[1] = 4/24 = P[2] = \dots = P[6]$ ,  $P[R \wedge 5] = 3/24$
- Note:  $P \geq 0$  always;  $P$  can be thought of in terms of proportions / frequency
- $P[\text{A particular ball}] = 1/24$

- R.V.s  $x := \# \text{ on ball } \in [6]$ ,  $y := \text{Colour on ball}$
- $P[x=1] = 4/24$ ,  $P[y=2] = 6/24$   $\{R=1, G=2, B=3\}$
- Suppose now they are not equally likely, but have prob. assigned  $\rightarrow$  No longer proportions  $y_{24} y_{24} y_{ab} y_{24} y_6 y_{ab}$   
 $y_{48} y_{24} y_{12} y_{ab} y_{48} y_{24}$   
 $y_{16} y_{24} y_{48} y_{ab} y_{12} y_{ab}$   
 $y_{ab} y_{24} y_{24} y_{48} y_{24} y_{12}$
- $P[x=1] = y_{24} + y_{48} + y_{16} + y_{ab} = 13/48$ ,  $P[y=2] = y_3$

## Rules of Probability

$\Omega$ : Sample Space,  $X$ : Any (discrete) R.V.

$S_x$ : Set of (numerical) values  $X$  could take  
 $\hookrightarrow$  called the support of the R.V.  $X$

$X(w)$ : Value of  $X$  on outcome  $w \in \Omega$

*Lazy notation*  $\xrightarrow{x \in S_x}$

$$P[x] = P[X=x] := \sum_{\substack{w \in \Omega \\ : X(w)=x}} p_w = \sum_{\substack{w \in \Omega \\ : X(w)=x}} P[w]$$

- If  $X$  is discrete valued, then this must hold:
  - $P[X=x] = 0 \Rightarrow x$  is an impossible value for  $X$
  - $P[X=x] = 1 \Rightarrow X$  almost surely takes  $\xrightarrow{\text{R.V. } X}$  the value  $x$
  - $\sum_{x \in S_x} P[X=x] = 1$   $\xrightarrow{x \text{ cannot remain undefined,}}$   
 $\leftarrow$  must take some valid value.  
 know this value but  $\exists$  some value  
 $\Rightarrow \forall x \in S_x, P[X=x] \leq 1$

## Probability Mass Function (PMF)

- $\hookrightarrow$  Tells us prob. of an RV taking a particular value
- For discrete RV  $X$ , its PMF  $f(\cdot)$  tells us for any  $x \in S_x$ ,  $f(x) = P[X=x]$
- $P_X[\cdot]$  used to emphasize that this PMF

is for  $x$  and not some other  $y$ .  $P(x)$  also used  
 Sampling from a PMF:  $x \sim P[x] / P_x$  or even  
 $x \sim P[x]$  means that we generated an outcome  
 $w \in \Omega$  acc to prob. dist. and are looking at  $x(w)$

Joint Probability  $P[x=1, y=2] = P[x=1 \wedge y=2]$   
 $\quad :=$  prob. that  $x=1 \wedge y=2$   
 $P[x=2 \wedge y=3] = 0$

PMF for joint Distribution  $\rightarrow P[x, y] = P_{x,y}[\cdot]$

★  $P[x=x] = \sum_{y \in S_y} P[x=x, y=y]$

★  $\because \sum_{x \in S_x} P[x=x] = 1 \Rightarrow \sum_{x \in S_x} \sum_{y \in S_y} P[x=x, y=y] = 1$

Can sample using PMF for JD too.

### Marginal Probability

When talking about only a subset of R.V.s, the PMF is called marginal prob. dist.

Given  $P_{x,y,z}$ , we can obtain

$$P_{x,z} = \sum_{y \in S_y} P_{x,y,z}; \quad P_x = \underbrace{\sum_{y \in S_y} \sum_{z \in S_z}}_{y \text{ & } z \text{ have}} P_{x,y,z}$$

→  $y$  has been marginalized out:  $\xrightarrow{\text{been marginalized out}}$

### Conditional Probability

$P[x=1 | y=2] :=$  Prob. of  $x=1$  among the samples in which  $y=2$   
 $\equiv \frac{\text{Proportion of samples w/ } y=2 \text{ & } x=1}{\text{Proportion of samples w/ } y=2}$

⇒

$$P[x=1 | y=2] = \frac{P[x=1 \wedge y=2]}{P[y=2]}$$

### PMF for Conditional Distribution

We had, for any  $x_0 \in S_x$ ,  $\sum_{y \in S_y} P[x=x_0, y=y] = P[x=x_0]$

$$\Rightarrow \sum_{y \in S_y} P[y=y | x=x_0] = \sum_{y \in S_y} \frac{P[y=y \wedge x=x_0]}{P[x=x_0]} = 1$$

- $P[y|x] = P_{y|x}[·|·]$ ; May sample  $y \sim P[y|x=x_0]$
- For sampling from  $P[y|x=x_0]$ , we consider the set of only those  $w \in \mathcal{N}$  where  $X(w)=x_0$  then sample an outcome  $w_0$  w/ prob.  $\frac{P[w_0]}{P[x=x_0]}$  & then return  $Y(w_0)$

## Marginal conditional Probability

$$P[x=x | z=z_0] = \sum_{y \in S_y} P[x=x, y=y | z=z_0]$$

$$P[y=y | x=x_0, z=z_0] = \frac{P[y=y \wedge (x=x_0 \wedge z=z_0)]}{P[x=x_0 \wedge z=z_0]}$$

$$= \frac{P[(y=y \wedge x=x) \wedge z=z]}{P[z=z]} \cdot \frac{P[z=z]}{P[x=x \wedge z=z]} = \frac{P[y=y, x=x | z=z]}{P[x=x | z=z]} =$$

- ★ Sum Rule  $P[x] = \sum_{y \in S_y} P[x, y]$   
(Marginalization Rule) aka Law of Total Probability
- ★ Product Rule  $P[x, y] = P[x|y] \cdot P[y]$   
(Conditioning Rule)  $P[x] = \sum_{y \in S_y} P[x, y] = \sum_{y \in S_y} P[x|y] \cdot P[y]$   
OR,  $P[x=x, y=y] = P[x=x | y=y] \cdot P[y=y]$
- ★ Chain Rule (Iterated Conditioning Rule)  
 $P[x, y, z] = P[x|y, z] \cdot P[y|z] \cdot P[z]$

## Bayes Theorem

$$P[x=x | y=y] = \frac{(P[x=x, y=y])}{P[y=y]}, \quad P[y=y | x=x] = \frac{(P[x=x, y=y])}{P[x=x]}$$

$$\Rightarrow P[x=x | y=y] \cdot P[y=y] = P[y=y | x=x] \cdot P[x=x]$$

$$\Rightarrow P[x=x | y=y] = P[y=y | x=x] \cdot P[x=x] / P[y=y]$$

- can create events using R.V.s too  
Ex:  $\{X=1\}$ ,  $\{Y=2\}$ ,  $\{X=1 \wedge Y=2\}$ ,  $\{X \leq x_0\}$ ,  $\{X < x_0\}$

Event Calculus → set of Rules

Let  $A, B$  be events,

$\neg A$  → negation/complement of  $A$  ( $A$  didn't happen)

$A \cup B$  → Union ( $A$  or  $B$  happened)

$A \cap B$  → Intersection ( $A$  &  $B$  happened)

De-Morgan's Laws (Always hold):

$$\neg(A \cup B) = \neg A \cap \neg B \quad \neg(A \cap B) = \neg A \cup \neg B$$

Complement Rule:  $P[\neg A] = 1 - P[A]$

Union Rule:  $P[A \cup B] = P[A] + P[B] - P[A \cap B]$

Intersection Rule:  $P[A \cap B] = P[A] + P[B] + P[\neg A \cap \neg B] - 1$

The above rules hold even for conditional probability:

$$P[\neg A | C] = 1 - P[A | C]$$

$$P[A \cup B | C] = P[A | C] + P[B | C] - P[A \cap B | C]$$

If  $C \Rightarrow A$  then  $P[A | C] = 1$ , and also

if  $C \Rightarrow \neg B$  then  $P[B | C] = 0$

Bayes' Theorem applies to events just as well.

Use indicator variables

\*  $\mathbb{I}\{E\} = 1$  if  $E$  is true else 0

## Independence of R.V.s

Two R.V.s  $X$  and  $Y$  are said to be independent

if  $\forall x \in S_x, y \in S_y$ ,

$$P[X=x, Y=y] = P[X=x] \cdot P[Y=y]$$

$$P[X=x | Y=y] = P[X=x]$$

$$P[Y=y | X=x] = P[Y=y]$$

If  $X$  and  $Y$  are independent, then we often write  $X \perp\!\!\!\perp Y$ . Similarly,  $X \not\perp\!\!\!\perp Y$

Conditional Independence

## Conditional Independence

$$X \perp\!\!\!\perp Y \mid Z \Leftrightarrow P[X=x, Y=y \mid Z=z] = P[X=x \mid Z=z] \cdot P[Y=y \mid Z=z]$$

- ★ If  $X, Y$  are independent, then it is not necessary that they continue to be independent even if conditioned on a 3rd R.V.
- ★ Even if  $X, Y$  are not independent, it's still possible that  $\exists$  3rd R.V.  $Z$  s.t.  $X \perp\!\!\!\perp Y \mid Z$ 
  - ↳ Commonly found in graphical models

**Random Vectors:** A collection of R.V.s arranged in an array  $X = [x_1, x_2, \dots, x_d]^T$

- Can be independent or correlated
- PMF/PDF of  $X$  = joint PMF/PDF of  $\{x_i\}_{i \in [d]}$
- Can talk about marginal / conditional prob. among  $x_1, \dots, x_d$ .

## Expectation of an R.V.

- Mean or the average value that R.V. takes and defined as

$$\mathbb{E}[x] = \sum_{x \in S_x} x \cdot P[x=x] = \mathbb{E} X$$

NOT most likely value

- ★ Linearity of Expectation  $\mathbb{E}[x+y] = \mathbb{E}[x] + \mathbb{E}[y]$

↳ Also the sum rule for expectation

- Scaling Rule:  $y = c \cdot x \Rightarrow \mathbb{E}[y] = c \cdot \mathbb{E}[x]$

$$* \mathbb{E}[x - \mathbb{E} x] = 0 = \mathbb{E} x - \mathbb{E} x \quad \because \mathbb{E} x \text{ is constant}$$

- \* Law of the Unconscious Statistician (LOTUS)

• We have R.V.  $x$  w/ PMF  $P_x$

& some  $g: S_x \rightarrow \mathbb{R}$  b  $y := g(x)$

$$\text{then } \mathbb{E} y = \mathbb{E} g(x) = \sum_{x \in S_x} g(x) \cdot P[x=x]$$

\* If R.V.s  $X$  and  $Y$  are independent, then we have  $\text{IE}[X \cdot Y] = \text{IE}[X] \cdot \text{IE}[Y]$

**Sample Mean:** Suppose R.V.  $X$  sampled again  $b$  agin, say  $n$  times. Values obtained, say  $x_1, x_2, \dots, x_n$ , we can estimate  $\text{IE}X$  if  $n$  sufficiently large

$\hat{\text{IE}}X = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

$\hat{\text{IE}}X = \arg \min_c \sum_{i=1}^n (x_i - c)^2 \rightarrow \text{Empirical}$

&

$\text{IE}X = \arg \min_c \text{IE}[(x - c)^2] \rightarrow \text{Theoretical}$

**Mode of an R.V.:** Value that R.V. takes w/ highest prob.

 $\text{mode}(x) := \arg \max_{x \in S_x} P[X=x]$ 

- An R.V. may have more than one mode
- Empirical mode:

$$\text{mode}(x_1, \dots, x_n) := \arg \max_{x \in S_x} \sum_{i=1}^n \mathbb{I}\{x_i = x\} = \arg \max_{x \in \{x_1, \dots, x_n\}} \sum_{i=1}^n \mathbb{I}\{x_i = x\}$$

**Median of an R.V.**

- A value  $m$  that satisfies  $P[X \geq m] \geq 0.5$  as well as  $P[X \leq m] \geq 0.5$
  - Empirically,  $(\# \text{ of samples } \geq m) = (\# \text{ of samples } \leq m)$
  - OR, arrange  $x_i$ 's  $x_1 \leq x_2 \dots \leq x_n$
- $n$  is odd  $\Rightarrow m = x_{(n+1)/2}$
- $n$  is even  $\Rightarrow m = (x_{n/2} + x_{1+n/2})/2$
- There may be infinitely many medians

\*  $m = \arg \min_c \sum_{i=1}^n |x_i - c|$

Empirical  $\leftarrow m = \arg \min_c \text{IE}[|X - c|] \rightarrow \text{Theoretical}$

**Variance:** How spread out are the values that an R.V. takes, how far from expectation

- For R.V.  $X$  with  $\mathbb{E}X = \mu$ , its variance is,  
 $\text{Var}[x] = \text{Var}[x] = \sigma^2 = \mathbb{E}[(x-\mu)^2] = \mathbb{E}[x^2] - (\mathbb{E}[x])^2$   
 $\Rightarrow \mathbb{E}[x^2] \geq (\mathbb{E}[x])^2, \quad \text{Var}[x] \geq 0$

Standard Deviation :=  $\sqrt{\text{Variance}} = \sigma$

Sample Variance:  $n$  samples  $x_1 \dots x_n$  of R.V.  $X$

$$\hat{\mu} = \frac{1}{n} \sum x_i, \quad M1: \hat{\sigma}^2 = \frac{1}{n} \sum (x_i - \hat{\mu})^2$$

$$M2: \hat{s} = \frac{1}{n} \sum x_i^2, \quad \hat{\sigma}^2 = \hat{s} - \hat{\mu}^2 \quad \hookrightarrow \text{Two passes}$$

Single pass / Running averages

Bad if  $\hat{s}$  &  $\hat{\mu}^2$  both very large & close  
{due to computers}

Covariance: Tells us how 2 R.V.s behave in tandem

- $\text{Cov}(x, y) := \mathbb{E}[(x - \mathbb{E}x)(y - \mathbb{E}y)] = \mathbb{E}[xy] - \mathbb{E}x \cdot \mathbb{E}y$
  - $\text{Cov}(x, x) = \text{Var}[x]$  and  $\text{Cov}(x, y) \leq 0$
  - Empirically,  $\hat{\mu}_x, \hat{\mu}_y$  precomputation
    - M1:  $\hat{\text{Cov}}(x, y) = \frac{1}{n} \sum (x_i - \hat{\mu}_x)(y_i - \hat{\mu}_y)$
    - M2:  $\hat{c} = \frac{1}{n} \sum x_i y_i$  and  $\hat{\text{Cov}}(x, y) = \hat{c} - \hat{\mu}_x \hat{\mu}_y$
- Same comments as before

## Rules of Variance

$b, c \in \mathbb{R}$  constants &  $x, y$  2 R.V.s

Constant Rule:  $\text{Var}[c] = 0$

Scaling Rule:  $\text{Var}[cx] = c^2 \cdot \text{Var}[x]$

Shift Rule:  $\text{Var}[c+x] = \text{Var}[x]$

Sum Rule:  $\text{Var}[x \pm y] = \text{Var}[x] + \text{Var}[y] \pm 2 \text{Cov}[x, y]$

## Rules of Covariance

Constant Rule:  $\text{Cov}[x, c] = 0$

Symmetry Rule:  $\text{Cov}[x, y] = \text{Cov}[y, x]$

Scaling Rule:  $\text{Cov}[bx, cy] = b \cdot c \cdot \text{Cov}[x, y]$

Shift Rule:  $\text{Cov}[x+b, y+c] = \text{Cov}[x, y]$

\* If  $x, y$  independent then

$\text{Cov}[x, y] = 0$



$$\Rightarrow \text{Var}[x+y] = \text{Var}[x] + \text{Var}[y]$$

**Correlation:** Normalized version of covariance

$$\rho_{x,y} = \text{Corr}[x,y] = \frac{\text{Cov}[x,y]}{\sqrt{\text{Var}[x] \cdot \text{Var}[y]}} \in [-1, 1]$$

$\rho_{x,y} = 0 \Rightarrow$  uncorrelated

Independent R.V.s  $\Rightarrow$  Uncorrelated

Empirically,

$$\hat{\rho}_{x,y} = \hat{\text{Cov}}(x,y) / \sqrt{\hat{\sigma}_x^2 \hat{\sigma}_y^2}$$

$\rho_{x,y} < 0 \rightarrow$  Typically,  $x$  takes larger value than its mean when  $y$  takes smaller value than its mean & vice-versa

$\rho_{x,y} > 0 \rightarrow \dots \text{larger} \dots \text{larger} \dots$

$\rho_{x,y} = 0 \rightarrow$  Typically, even if  $x$  takes a value larger than its own mean,  $y$  may take smaller or larger values than its mean

## Conditional Formatting

- [.] used to express how one quantity behaves when some other quantities are fixed
- Other quantities could be R.V.s themselves or constants.

Ex:  $P$  of misclassifying  $(x, y) \sim D$  if model  $w$  used, i.e.,  $P[y \cdot w^T x < 0 | w]$

Conditional Expectation  $E[x | y=y_0] := \sum_{x \in S_x} x \cdot P[x=x | y=y_0]$

Conditional Variance  $\text{Var}[x | y=y_0] := E[(x-\mu)^2 | y=y_0]$

Conditional Covariance

$$\mu = E[x | y=y_0]$$

$$\text{Cov}[x, y | z=z_0] = E[(x-\mu_x)(y-\mu_y) | z=z_0]$$

$$= E[xy | z=z_0] - \mu_x \mu_y, \quad \mu_x = E[x | z=z_0]$$

Conditional Mode  $\text{mode}[x | y=y_0] = \arg \max_{x \in S_x} P[x=x | y=y_0]$   
These defns don't require independence.

Rules of expectation hold if all expressions are systematically conditioned

$$\mathbb{E}[X+Y|Z=z_0] = \mathbb{E}[X|Z=z_0] + \mathbb{E}[Y|Z=z_0],$$

$$\mathbb{E}[c \cdot X|Z=z_0] = c \cdot \mathbb{E}[X|Z=z_0],$$

$$\mathbb{E}[g(x)|Z=z_0] = \sum_{x \in S_x} g(x) \cdot P[X=x|Z=z_0],$$

$$\text{If } X \perp\!\!\!\perp Y | Z \text{ then } \mathbb{E}[X \cdot Y|Z=z_0] = \mathbb{E}[X|Z=z_0] \cdot \mathbb{E}[Y|Z=z_0]$$

## Expectation of an R.V.

$$\mathbb{E} X = [\mathbb{E} X_1, \mathbb{E} X_2, \dots, \mathbb{E} X_d]^T$$

$$\mathbb{E}[X+Y] = \mathbb{E} X + \mathbb{E} Y, \quad \mathbb{E}[c \cdot X] = c \cdot \mathbb{E}[X], \quad c \in \mathbb{R} \text{ constant}$$

$$a \in \mathbb{R}^d \text{ constant}, \quad \mathbb{E}[a^T X] = a^T \mathbb{E}[X]$$

$$A \in \mathbb{R}^{n \times d} \text{ constant}, \quad \mathbb{E}[AX] = A \mathbb{E}[X]$$

Mode easy to define. Median not so easy:

$$\text{Definition 1: } \text{med}(x) = [\text{med}(x_1) \dots \text{med}(x_d)]^T$$

D2: minimizer of absolute distance (L1 norm)

$$\text{med}(x) = \arg \min_{v \in \mathbb{R}^d} \mathbb{E}[\|x - v\|_1]$$

$$\text{Cov}(x) = \begin{bmatrix} \mathbb{V}x_1 & \text{Cov}(x_1, x_2) & \dots & \text{Cov}(x_1, x_d) \\ \text{Cov}(x_2, x_1) & \mathbb{V}x_2 & & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \text{Cov}(x_d, x_1) & \dots & \dots & \mathbb{V}x_d \end{bmatrix} = \mathbb{E}[(x-\mu)(x-\mu)^T] = \mathbb{E}[xx^T] - \mu\mu^T, \quad \mu = \mathbb{E}[x]$$

$$\text{Cov}(cx) = c^2 \text{Cov}(x), \quad x \text{ is a vector of R.V.s}$$

$$\text{If } x \in \mathbb{R}^m, y \in \mathbb{R}^n \text{ be two R. vectors, then } \mathbb{E}[\text{Cov}(x, y)] = \mathbb{E}[(x-\mu_x)(y-\mu_y)^T] = \mathbb{E}[xy^T] - \mu_x\mu_y^T,$$

$$a \in \mathbb{R}^d \text{ constant, then } \mathbb{E}[\text{Cov}(a^T x)] = a^T \text{Cov}(x) a$$

$$A \in \mathbb{R}^{n \times d} \text{ constant, then } \mathbb{E}[\text{Cov}(Ax)] = A \text{Cov}(x) A^T \in \mathbb{R}^{n \times n}$$

## Continuous R.V.s

- R.V. that can take infinitely many possible values that are not discrete but cts., i.e., support is  $\mathbb{R}$  or some subset of  $\mathbb{R}$

- Instead of PMF, use PDF (prob. dist. func.)

- $f = \text{PDF}: S_x \rightarrow \mathbb{R}_+$  + can. take values  $> 1$

as well but must not be -ve

- For any  $x \in S_x$ ,  $f_x(x)$  tells us how likely is  $X$  to take a value around  $x$ , ie, for some  $\delta > 0$ , we have  $P[X \in [x-\delta, x+\delta]] \approx f_x(x) \cdot 2\delta$   
 $\int_{x-\delta}^{x+\delta} f_x(t) dt$  : Exact formula

More generally, we have,

$$P[X \in [a, b] \subseteq S_x] = \int_a^b f_x(t) dt$$

- PDF  $f$  of an R.V.  $X$  satisfies

$$f_x(x) \geq 0 \quad \forall x \in S_x \quad \& \quad \int_{S_x} f_x(t) dt = 1$$

Expectation :  $E[X] := \int_{S_x} t \cdot f_x(t) dt$

LOTUS :  $E[g(X)] = \int_{S_x} g(t) f_x(t) dt$

Variance :  $V[X] := \int_{S_x} (t - E[X])^2 f_x(t) dt = \int_{S_x} t^2 f_x(t) dt - (E[X])^2$

Joint PDFs :  $f_{x,y}$ ,  $[a, b] \subseteq S_x$ ,  $[c, d] \subseteq S_y$

$$P[X \in [a, b], Y \in [c, d]] = \int_a^b \int_c^d f_{x,y}(s, t) ds dt$$

→ Make sense even if one cts. & other discrete

Marginal PDF :  $f_x(x) = \int_{S_y} f_{x,y}(x, t) dt$

conditional probabilities also make sense

$P[X \in [a, b] | Y=y]$  can be found using limits or  
in general zero ← Radon-Nikodym derivative

conditional expectations, (co)variances also defined

All rules of Probability, Expectation, (Co)Variance  
continue to hold

## Probability Distributions

### Bernoulli Distribution

- Support :  $\{0, 1\} \rightsquigarrow$  Useful in binary classification
- PMF of R.V.  $Y$  uniquely specified by  
 $P[Y=1] = p \Leftrightarrow P[Y=0] = 1-p$

Bias / Success probability

- Mean :  $p$ , Mode : 1 if  $p > 0.5$ , 0 if  $p < 0.5$ ,  $\frac{1}{2}$  if  $p = 0.5$

$$\text{Variance: } p(1-p)$$

## Rademacher Distribution

Support:  $\{-1, 1\}$

$X$  distributed over Bern  $\Leftrightarrow 2X-1$  over Rade

$Y$  dist. over Rade  $\Leftrightarrow (Y+1)/2$  over Bern

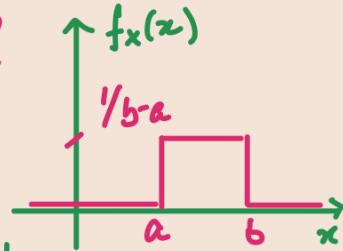
$P[Y=1] = p$ , Mean =  $2p-1$ , Mode: 1 if  $p > 0.5$ ,

Variance:  $4 \cdot p \cdot (1-p)$  -1 if  $p < 0.5$ ,  $\{-1, 1\}$   $p=0.5$

## Uniform Distribution

$X$  cts. R.V. w/  $S_x = [a, b] \subseteq \mathbb{R}$

PDF  $f_x(x) = \begin{cases} \frac{1}{b-a}, & x \in [a, b] \\ 0, & \text{otherwise} \end{cases}$



$f_x(x) \geq 0$  if  $x \in S_x$  &  $\int_{S_x} f(t) dt = 1$

$E[X] = (a+b)/2$ ;  $V[X] = (b-a)^2/12$  ↑ as  $(b-a) \uparrow$

$\text{UNIF}([a, b])$  used to denote uniform dist.

## Gaussian (aka Normal) Distributions

$S_x = \mathbb{R}$

$$f_x[x | \mu, \sigma^2] = \frac{\frac{-(x-\mu)^2}{2\sigma^2}}{\sqrt{2\pi\sigma^2}}$$

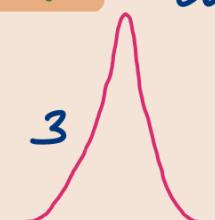
Takes two numbers to describe

$\mu \in \mathbb{R}$ ;  $\sigma^2$ : non-negative real number

Often written as

$$\mathcal{N}_x(x; \mu, \sigma^2) \quad \text{or} \quad \mathcal{N}(x; \mu, \sigma^2)$$

Mean =  $E[X] = \mu$  = Mode = Median,  $V[X] = \sigma^2$  Conditioned to constants



$$Z = cX \rightarrow \mathcal{N}_z(\cdot; c\mu_x, c^2\sigma_x^2)$$

$$W = X + Y \rightarrow \mathcal{N}_w(\cdot; \mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$$

$$V = X + Y \rightarrow \mathcal{N}_v(\cdot; \mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$$

$X, Y$  Gaussian  
 $\delta$  independent

$\sqrt{1 - x^2} \leq \sqrt{1 - (\mu_x + t\sigma_x)^2} \leq \sqrt{1 - \mu_x^2 - 2t\mu_x\sigma_x + t^2\sigma_x^2}$

$$\Pr[X \geq \mu_x + t\sigma_x] = \Pr[X \leq \mu_x - t\sigma_x] \leq e^{-t^2/2}$$

As  $\sigma_x \downarrow$ , R.V. gets more concentrated around its mean

★ **68-95-99.7 rule**

$$\Pr[|\mathbf{x} - \mu| \leq \sigma_x] \approx 0.68, [\leq 2\sigma_x] \approx 0.95, [\leq 3\sigma_x] \approx 0.997$$

## Gaussian Random Vector

- Need only the mean  $\mu \in \mathbb{R}^d$  and the covariance  $\Sigma \in \mathbb{R}^{d \times d}$  to be specified  $N(\mu, \Sigma)$
- $$\Pr[\mathbf{x} | \mu, \Sigma] = \frac{e^{-\frac{1}{2}(\mathbf{x}-\mu)^\top \Sigma^{-1}(\mathbf{x}-\mu)}}{\sqrt{(2\pi)^d |\Sigma|}}$$

\* Special case  $\mu=0, \Sigma=I_d$  called standard Gaussian / Normal dist,

$$\Pr[\mathbf{x} | 0, I_d] = \frac{e^{-\frac{1}{2}\|\mathbf{x}\|_2^2}}{\sqrt{(2\pi)^d}} = \prod_{i=1}^d \frac{e^{-\frac{x_i^2}{2}}}{\sqrt{2\pi}}$$

However,  $e^{-\frac{x_i^2}{2}/\sqrt{2\pi}}$  is simply  $N(0, 1)$ , i.e,

$$\Pr[x_1, \dots, x_d | 0, I] = \prod \Pr[x_i | 0, 1]$$

⇒ All  $d$  co-ordinates of a standard Gaussian R.V. are independent

- Given a Gaussian vector  $\mathbb{R}^d \ni \mathbf{x} \sim N(\mu, \Sigma)$ :
  - Every coordinate of  $\mathbf{x}$  is a Gaussian R.V.
  - Need not be independent
  - Holds true even if conditioned on all other co-ords.
- Consider  $j \in [d]$ 
  - $x_j$  distributed as  $N(\mu_j, \Sigma_{jj})$
  - Given  $x_k = v_k + k \neq j$ ,  $x_j$  still Gaussian, expression
- If  $a \in \mathbb{R}^d$  const,  $\mathbb{R} \ni a^\top \mathbf{x} \sim N(a^\top \mu, a^\top \Sigma a)$  complicated
- If  $A \in \mathbb{R}^{n \times d}$  const,  $\mathbb{R}^n \ni Ax \sim N(A\mu, A\Sigma A^\top)$

## Laplacian Distribution

- Concentrate much more strongly

around mean than Gaussian

$$f_x[x | \mu, \sigma] = \frac{e^{-|x-\mu|/\sigma}}{2\sigma}$$



Mean = Mode = Median =  $\mu$ . Variance =  $2\sigma^2$

$y = ax + b$  also Laplacian w/  $\mu_y = a \cdot \mu_x + b$ ,  $\sigma_y = a \cdot \sigma_x$



## Uncertainty

Start State, Goal States, Good Ques, Bad Ques.

•	•	•	•
•	•	•	•
•	•	•	•
•	•	•	•

•	•	•	•
•	•	•	•
•	•	•	•
•	•	•	•

•	•	•	•
•	•	•	•
•	•	•	•
•	•	•	•

•	•	•	•
•	•	•	•
•	•	•	•
•	•	•	•

- Asking trivially true questions : Useless
- Asking rarely true questions : Useless
- Must ask questions that eliminate 'wrong' answers quickly

**Entropy ( $H$ ):** If we have a set of  $n$  words, then that set has an entropy of  $\log_2 n$

- Larger sets have larger entropy ; 1 word : 0 entropy
- More generally, if there is a set  $S$  of  $n$  elements of  $C$  types with  $n_c$  elements of type  $c$ , then its entropy :

$$H(S) := - \sum_{c \in [C]} \frac{n_c}{n} \log_2 \left( \frac{n_c}{n} \right) = - \sum_{c \in [C]} p_c \log_2 p_c$$

where  $p_c$  is the proportion of elements of type  $c$

- \* A pure set, ex:  $p = [0 0 \dots 1 \dots 0 0]$  has 0 entropy whereas a set with same number of elements of each class i.e.,  $p = [\frac{1}{C} \frac{1}{C} \dots \frac{1}{C}]$ :  $\log_2 C$  entropy

What is a good question? — Depends on the application

- ID3 (iterative dichotomizer 3)

↳ reduces confusion/entropy the most

- Suppose question splits  $S$  into  $K$  subsets  $S_1, \dots, S_K$ ,  
 $\left\{ \begin{array}{l} S_i \cap S_j = \emptyset \quad \forall i \neq j, \\ \bigcup_{k \in [K]} S_k = S, \\ n_k := |S_k|, \quad \sum_{k \in [K]} n_k = n := |S| \end{array} \right.$

Then the entropy of  
the collection of sets is:

$$\sum_{k \in [K]} \frac{n_k}{n} \cdot H(S_k)$$



- Can be interpreted as 'average' or 'weighted' entropy

Ex:  $\{4096 \text{ words}\} \xrightarrow[\text{split}]{\text{Question}} \{1024, 1024, 1024, 1024\}$

Old entropy:  $\log_2 4096 = 12$

New entropy:  $\frac{1024}{4096} \log_2 1024 + \dots = 4 \cdot \frac{1024}{4096} \log_2 1024 = 10$   
We gained  $12 - 10 = 2$  bits of information

- \* Information gained this way can be added up.

Can be proved mathematically that the chosen definition of entropy is the only one satisfying 3 intuitive requirements. Suppose an event occurs with prob.  $p$  and we wish to measure the information from that event's occurrence, say  $I(p)$ , s.t.:

1) Rare event conveys no information:

2) More common the event,

less info it conveys:

$$p \uparrow \Rightarrow I(p) \downarrow$$

3) Info conveyed by two indep. events adds up:

$$I(p_1 \cdot p_2) = I(p_1) + I(p_2)$$

$$\Rightarrow I(p) = -\log_b p \Rightarrow H(S) = \sum_{c \in [C]} p_c \cdot I(p_c)$$

$b=2$ : bits (binary digits),  $b=e$ : nits (natural digits)

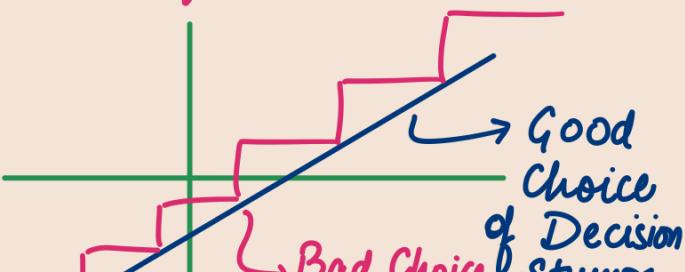
$b=10$ : digits (decimal digits / hartleys)

## ID3 Algorithm

- With  $T$  as set of all train points, create a root node  $r$ , and call  $\text{train}(r, T)$
- $\text{train}(\text{node } n, \text{set } S)$ 
  - If  $S$  is sufficiently pure or small, make  $n$  a leaf, decide a simple leaf action (e.g. most popular class, label popularity vector, etc.) and return
  - Else, out of available choices, choose the splitting criteria (e.g. a single feature) that causes maximum information gain / reduces entropy the most
  - Split  $S$  along that criteria to get  $S_1, \dots, S_k$  partition of  $S$  (e.g. if that feature takes  $k$  distinct values)
  - Create  $K$  child nodes  $c_k, k \in [K]$  and call  $\text{train}(c_k, S_k) \forall k \in [K]$
- ID3 does not ensure a balanced tree but the balance is decent.
- Several augmentations to this algorithm: C4.5, C5.0 allow real-valued features, missing features, boosting, etc.

## Careful Use of DTs

- Can be tweaked to give very high training accuracies — can badly overfit if tree too large
- Choice of decision stumps is critical:
  - PUF: single linear model DT will struggle and



eventually overfit



Decision Stumps

Probabilistic ML techniques, given a data point, do not output a single label, they instead output a distribution over all possible labels

- For binary classification, output a PMF over  $\{-1, 1\}$
- For multiclassification, output a PMF over  $\{1, 2, \dots, C\}$
- For regression, output a PDF over  $\mathbb{R}$
- The probability of a label indicates how likely does the ML model think that label is the correct one for that data point.
- The algo. is allowed to output a possibly different PMF / PDF for every data point. However, the support is always the set of all possible labels (even very unlikely labels)

## Prob. ML for Classification

- Say we have learnt a PML model  $w$  which, for a data point  $x$ , gives us a PMF  $P[Y=y|x, w]$  over the set of all possible labels, say  $Y = \{-1, +1\}$  for bin,  $Y = [C]$  for multiclassification
- May use the mode if single label predicted:  $\hat{y} = \arg \max_{y \in Y} P[Y=y|x, w]$
- May use the median/ mean as well — Bayesian ML exploits this
- Use  $P[Y=\hat{y}|x, w]$  to find out if the ML model is confident about its prediction or confused.  
{Confidence  $\nRightarrow$  Correctness}
- In many ML applications (e.g. active learning), if we find that the model is unsure, we can switch to another model or ask a human to step in

- May use variance of  $P[y|x, w]$  as well  
(low variance: very confident; high-var: less confident)

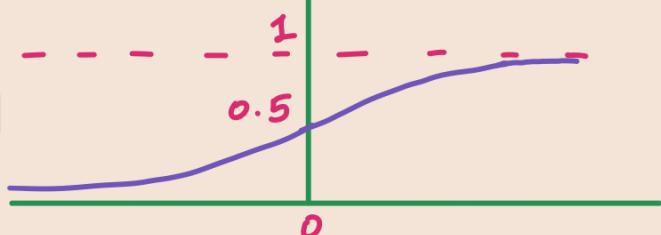
## Prob. Binary Class.

- Find a way to map every data point  $x$  to a Rademacher distribution, i.e., map each  $x$  to a prob  $p_x \in [0, 1]$ 
  - Gives us a PMF  $[1-p_x, p_x]$
- If using mode predictor,  $\hat{y}$ , then this PMF will give us the correct label only if:
  - When true label of  $x$  is +1,  $p_x > 1-p_x \Leftrightarrow p_x > 0.5$
  - " " " " " " " " -1,  $p_x < 1-p_x \Leftrightarrow p_x < 0.5$
  - $p_x = 0.5 \Rightarrow$  model is completely confused  
→  $x$  is on the decision boundary
- As usual, we want a healthy margin:  
True label of  $x$  is +1, we want  $p_x \gg 0.5 \Rightarrow p_x \approx 1$   
" " " " " " " " -1, " " " " " " " "  $p_x \ll 0.5 \Rightarrow p_x \approx 0$

Mapping feature vectors of  $x$  to  $p_x \in [0, 1]$

- could treat it as a regression problem
- $p_x \in \mathbb{R} \rightsquigarrow$  Will need to modify the training set (change all -1 labels to 0)
- could use DT
- Using linear models presents a challenge
  - If we learn a linear model  $w$  using regression,  $w^T x < 0 / w^T x > 1$  may happen
- DT won't run into this problem  $\because p_x$  always in  $[0, 1] \quad \{ \text{Average of } 0s \& 1s \in [0, 1] \text{ too} \}$
- can use logistic regression in linear models for PML

$$\sigma(t) = \frac{1}{1+e^{-t}} = \frac{e^t}{1+e^t} = 1 - \sigma(-t)$$



- Learn a linear model  $w$  (possibly with a hidden/explicit bias) & map  $x \rightarrow \sigma(w^T x)$
- Valid PMF  $\sigma(t) > 0.5 \Leftrightarrow t > 0$ ;  $\sigma(t) < 0.5 \Leftrightarrow t < 0$   
 $\sigma(t \rightarrow \infty) \rightarrow 1$ ;  $\sigma(t \rightarrow -\infty) \rightarrow 0$   
 $\Rightarrow$  Sigmoidal map will predict  $p_x \approx 1$  if  $w^T x \gg 0$   
and  $p_x \approx 0$  if  $w^T x \ll 0$
- Several other such wrapper/quashing/link/activation functions like tanh, ramp, ReLU
- How to learn such a  $w$ ?

## Likelihood

- Suppose we have a model  $w$  (assume hidden b)
- Given a data point  $(x^t, y^t)$ ,  $x^t \in \mathbb{R}^d$ ,  $y^t \in \{-1, +1\}$ ,  
the use of a sigmoidal map gives us a Rademacher PMF  $P[y | x^t, w]$
- The probability that this PMF gives to the correct label, i.e.,  $P[y^t | x^t, w]$  is called the "likelihood" of this model to this data point.
- Easy to show that  $P[y^t | x^t, w] = \sigma(y^t \cdot w^T x^t)$
- If we have several data points  $(x^1, y^1), \dots, (x^n, y^n)$  then we define the likelihood of  $w$  w.r.t. entire dataset as  $P[y^1, \dots, y^n | x^1, \dots, x^n, w]$
- Usually assumed that the data points are independent:  
 $P[y^1, \dots, y^n | x^1, \dots, x^n, w] = \prod_{i=1}^n P[y^i | x^i, w] = \prod_{i=1}^n \sigma(y^i \cdot w^T x^i)$
- Not always independent
- The expression  $P[y^i | x^i, w]$  tells us if the model  $w$  thinks the label  $y^i$  is a very likely label given the feature vector  $x^i$ .

label given the feature vector  $x$

Similarly  $P[y^1 \dots y^n | x^1 \dots x^n w]$

If we trust the training data then we should learn a model  $w$  that considers training labels to be very likely.

## Maximum Likelihood Estimator (MLE)

The model that gives highest likelihood to observe labels

$$\hat{w}_{MLE} = \arg \max_{w \in \mathbb{R}^d} \prod_{i=1}^n P[y^i | x^i, w]$$

If we're learning the model using a sigmoidal map, then product of several such values ( $[0, 1]$ ) can be extremely small.

→ Take log & exploit  $\max f(w) = \max_w \ln(f(w))$

Working w/ products can be numerically unstable

$$\begin{aligned}\hat{w}_{MLE} &= \arg \max_{w \in \mathbb{R}^d} \ln \left( \prod_{i=1}^n \sigma(y^i \cdot w^T x^i) \right) \\ &= \arg \min_{w \in \mathbb{R}^d} \sum_{i=1}^n \underbrace{\ln(1 + e^{-y^i \cdot w^T x^i})}_{\text{Also called negative log-likelihood}}\end{aligned}$$

\* Logistic loss function pops out automatically when learning a model that maximizes likelihood

## Prob. Multiclassification

Suppose  $C$  classes, then for every data point, we would like to output a PMF over the support  $[C]$

→ Assign a +ve score to all classes & normalize

To convert any score to +ve — exponentiate

learn  $C$  models  $w^1, \dots, w^C$ , given a point  $(x^t, y^t)$ .  
 $x^t \in \mathbb{R}^d$ ,  $y^t \in [C]$ . assign a +ve score per class  
 $\eta_c = e^{w^c \cdot x^t}$ ; then normalize to obtain the PMF  $P[y^t | x^t, \{w^c\}] = \eta_{y^t} / \sum \eta_c$  for any  $y \in [C]$

Likelihood:  $P[y^t | x^t, \{w_c\}] = \eta_{y^t} / \sum \eta_c$

Log-likelihood:  $\ln(\eta_{y^t} / \sum \eta_c)$

Distributions for  $C > 2$  called Multinoulli / Categorical distributions

## Softmax Regression

- If we now want to learn the MLE, we would have to find

$$\{\hat{w}_1^{\text{MLE}}, \dots, \hat{w}_C^{\text{MLE}}\} = \arg \max_{w^1, \dots, w^C \in \mathbb{R}^d} \prod_{i=1}^n P[y^i | x^i, w]$$
$$= \arg \max_{w^1, \dots, w^C \in \mathbb{R}^d} \prod_{i=1}^n \frac{\eta_{y^i}^i}{\sum_{c=1}^C \eta_c^i}, \quad \eta_c^i = e^{w_c^T x^i}$$
$$= \arg \min_{w^1, \dots, w^C \in \mathbb{R}^d} \sum_{i=1}^n -\ln \left( \frac{\eta_{y^i}^i}{\sum_{c=1}^C \eta_c^i} \right)$$

Cross-entropy / softmax loss function

- But 3 other ways to do prob. multiclass. as well

## General Recipe for MLE Algorithms

- Label set  $y$ , find a way to map data features  $x$  to PMFs  $P[\cdot | x, m]$  w/ support  $y$  parameters of the model

$P[\cdot | x, m] \rightarrow$  Likelihood function

$-\ln P[\cdot | x, m] \rightarrow$  Log-likelihood function

- Given data  $\{(x^i, y^i)\}_{i=1}^n$ , find model parameters that maximize likelihood function, ie, think that the training labels are very likely

$$\hat{m}^{\text{MLE}} = \arg \min_m \sum_{i=1}^n -\ln P[y^i | x^i, m]$$