# CS 189/289

Today's lecture:

- Maximum likelihood estimation (MLE)
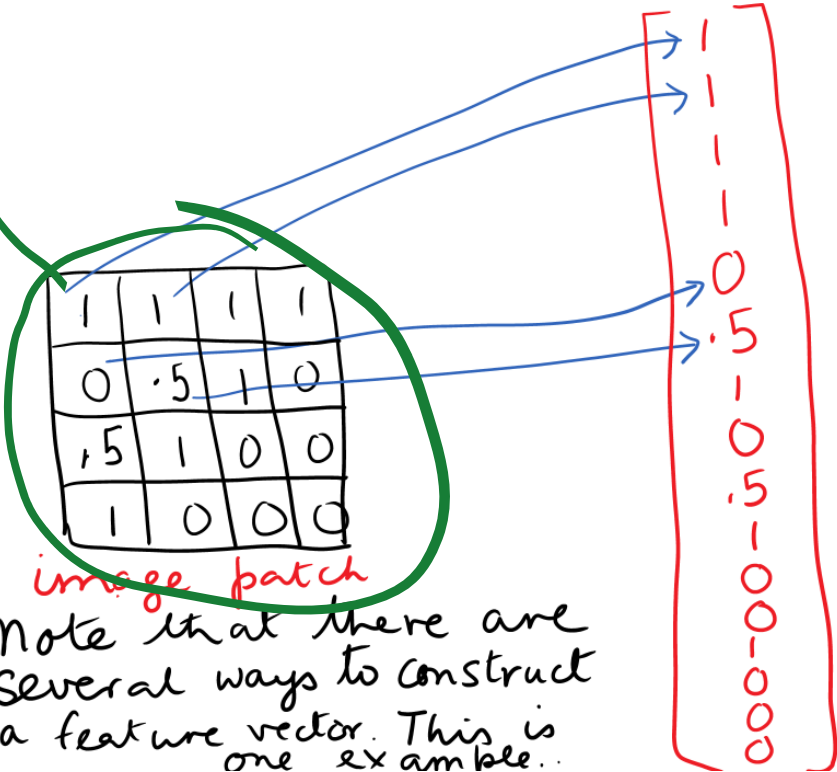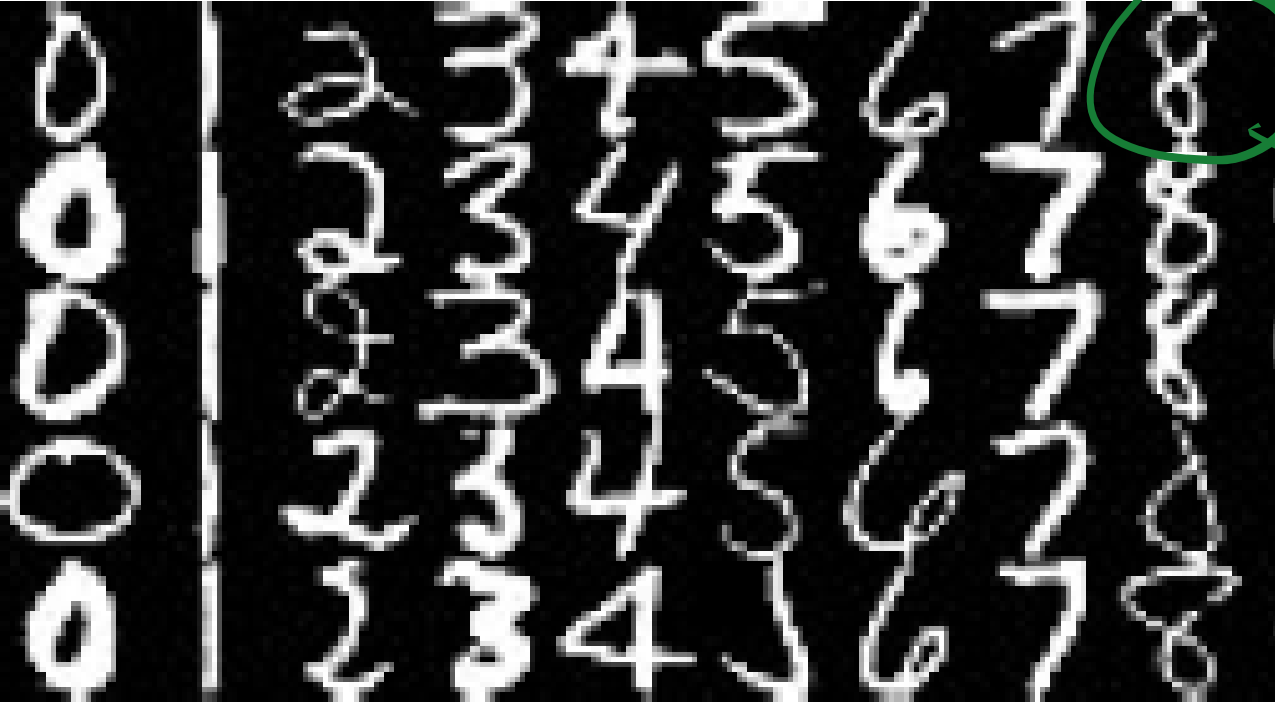
# Recall from last class:

Problem of digit classification from handwriting: is  a "7", yes or no?



- 60K training examples of digits (6K per class)
- Each digit is a 28 x 28 pixel grey level image.

# Recall from last class:

Problem of digit classification from handwriting: is  a "7", yes or no?
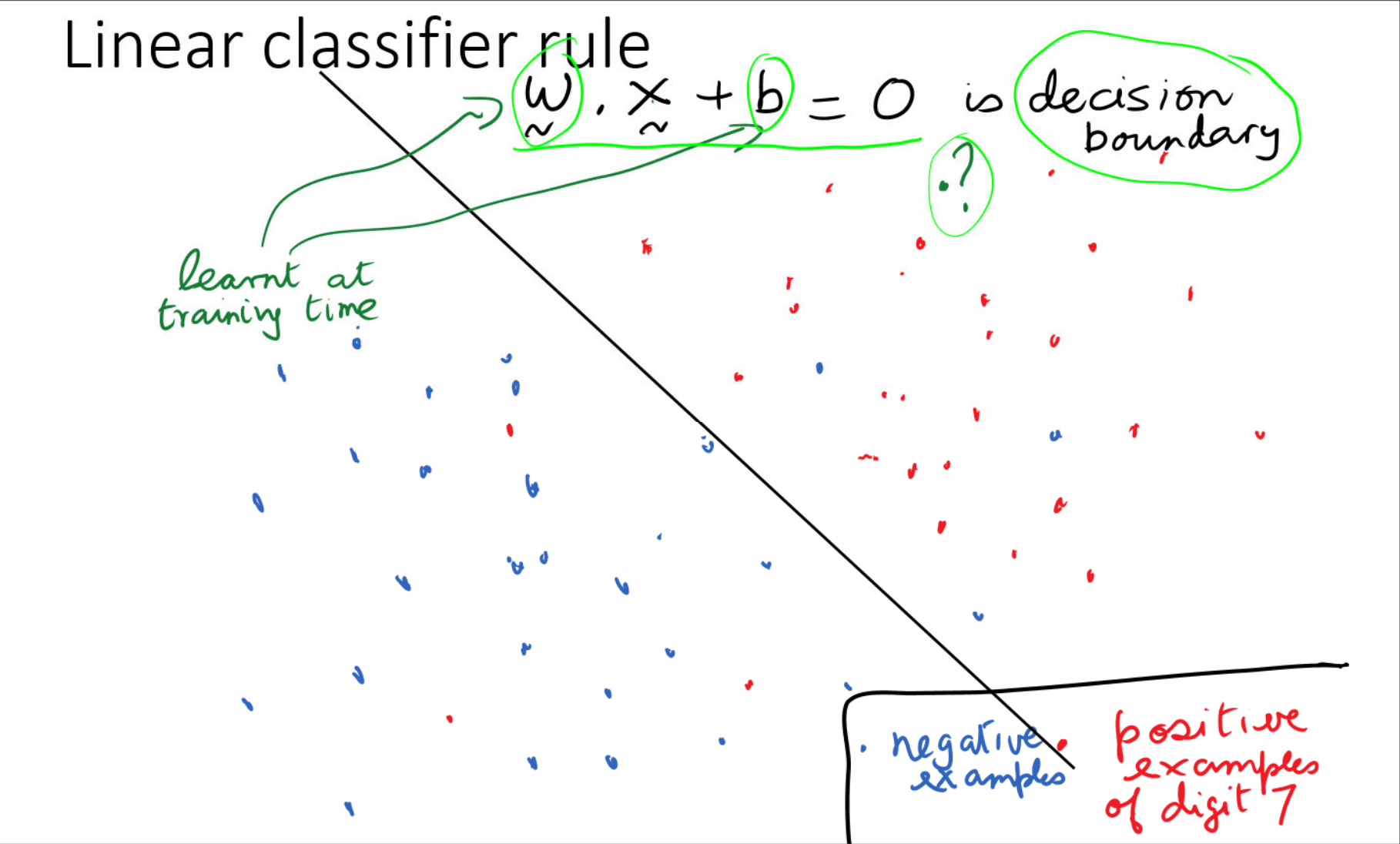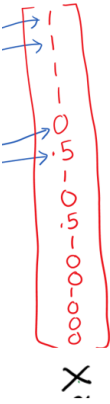


image patch

note that there are several ways to construct a feature vector. This is one example.
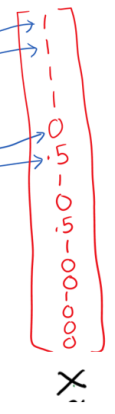
Feature vector $\mathbb{R}^{16}$

- 60K training examples of digits (6K per class)
- Each digit is a 28 x 28 pixel grey level image.

# *Recall from last class:*



Linear classifier rule

$$\underset{\sim}{w} \cdot \underset{\sim}{x} + b = 0 \quad \text{is} \quad \text{decision boundary}$$

learnt at training time

negative examples

positive examples of digit 7

# *Recall from last class:*



Linear classifier rule

$\underset{\sim}{w} \cdot \underset{\sim}{x} + b = 0$ is decision boundary
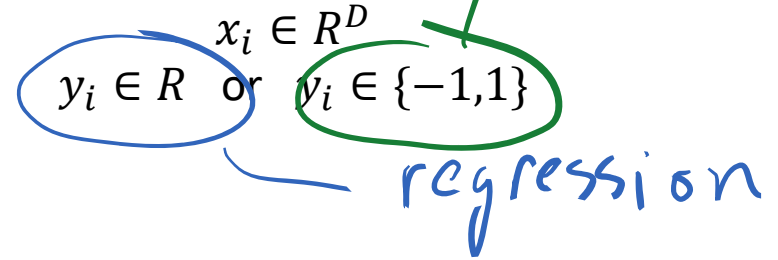
.?

- One of the main ways to "learn" (aka estimate) the setting of "good" parameters in statistical models:
- Principle of *Maximum Likelihood Estimation* (MLE).

negative examples

positive examples of digit 7

# ML: main concepts

- Training data set: $D = \{(x_i, y_i)\}_{i=1}^N$

$x_i \in R^D$

$y_i \in R$ or $y_i \in \{-1,1\}$

*classification*

*regression*

# ML: main abstract ideas

- Training data set:

$$D = \{(x_i, y_i)\}_{i=1}^{N}$$

$x_i \in R^D$
$y_i \in R$ or $y_i \in \{-1,1\}$

SUPERVISED

$$D = \{(x_i)\}_{i=1}^{N}$$

$x_i \in R^D$

UNSUPERVISED

"label"
provides
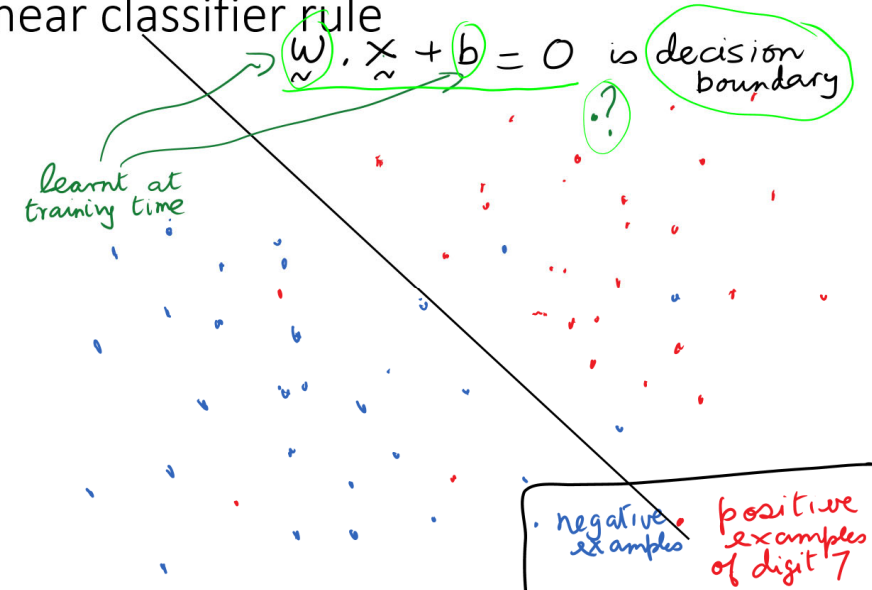supervision

# ML: main abstract ideas

- Training data set:

$$D = \{(x_i, y_i)\}_{i=1}^{N}$$

$$x_i \in R^D$$
$$y_i \in R \quad \text{or} \quad y_i \in \{-1,1\}$$

- Model class:
aka hypothesis class

$$f(x|w, b) = w^T x + b$$

**Linear Models**

Linear classifier rule

$\underset{\sim}{w} \cdot \underset{\sim}{x} + b = 0$ is decision boundary

learnt at training time

negative examples

positive examples of digit 7

# ML: main abstract ideas

- Training data set:

$$D = \{(x_i, y_i)\}_{i=1}^N$$

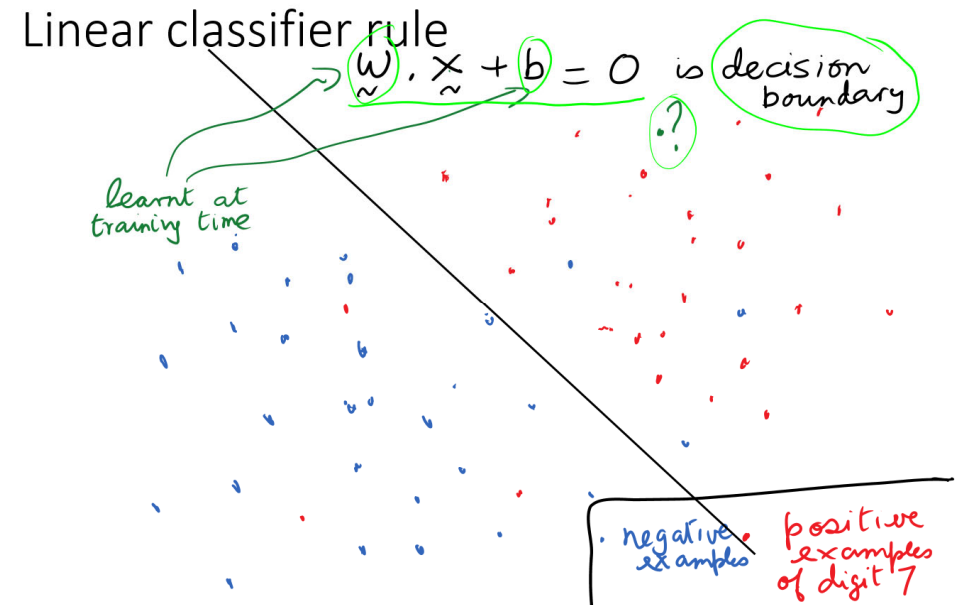$$x_i \in R^D$$
$$y_i \in R \text{ or } y_i \in \{-1,1\}$$

- Model class:
aka hypothesis class

$$f(x|w, b) = w^T x + b$$

**Linear Models**

Linear classifier rule

$\underset{\sim}{w} \cdot \underset{\sim}{x} + b = 0$ is decision boundary

learnt at training time

?

- **Optimization goal**: find "good" values of parameters $(w, b)$.

But was does "good" mean?

negative examples

positive examples of digit 7

# ML: main abstract ideas

- Training data set:

$$D = \{(x_i, y_i)\}_{i=1}^{N} \qquad \begin{matrix} x_i \in R^D \\ y_i \in R \ \text{ or } \ y_i \in \{-1,1\} \end{matrix}$$

- Model class:
aka hypothesis class

$$f(x|w,b) = w^T x + b$$

**Linear Models**

- Loss Function:

$$L(a,b) = (a-b)^2$$

**Squared Loss**

- Learning Objective:

$$\underset{w,b}{\operatorname{argmin}} \sum_{i=1}^{N} L\big(y_i, f(x_i \mid w,b)\big)$$

Optimization Problem

# Maximum Likelihood Estimation (MLE)

Linear classifier rule

$\underline{w} \cdot \underset{\sim}{x} + b = 0$ is decision boundary

learnt at training time

negative examples

positive examples of digit 7

This principle gives a useful, principled and widely-used loss function to estimate parameters of statistical models (from linear regression, to neural networks, and beyond).

- Training data set:
  $$D = \{(x_i, y_i)\}_{i=1}^N \qquad \begin{array}{c} x_i \in R^D \\ y_i \in R \quad \text{or} \quad y_i \in \{-1,1\} \end{array}$$

- Model class:
  aka hypothesis class
  $$f(x|w,b) = w^T x + b \qquad \text{Linear Models}$$

- Loss Function:
  $$L(a,b) = (a-b)^2 \qquad \text{Squared Loss}$$

- Learning Objective:
  $$\underset{w,b}{\text{argmin}} \sum_{i=1}^N L(y_i, f(x_i \mid w, b))$$

  Optimization Problem
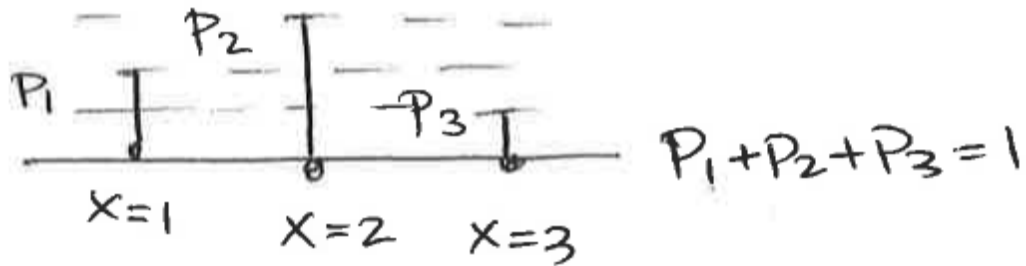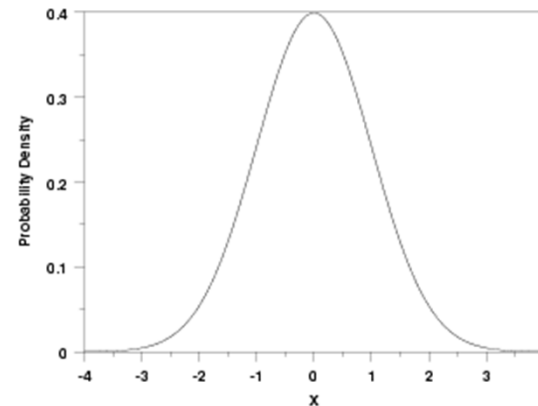
RVs!

# Reminder: probability distributions

Random variable (RV) is a function: $x \to \mathbb{R}$   e.g. $p(\text{heads}) = 0.5$

1. Discrete RV, e.g. coin toss heads/tails.
2. Continuous RV, e.g. height

Discrete RVs have a Probability Mass Function (PMF)

Continuous RVs have a Probability Density Function (PDF)



$P_1 + P_2 + P_3 = 1$

integrates to 1

# e.g. distributions of discrete RVs

1. <u>Bernouilli RV</u>—model the toss of a coin that can be biased
   $P(heads) = p, \quad P(tails) = 1 - p$, parameter is $p$.

# e.g. distributions of discrete RVs

1. <u>Bernouilli RV</u>—model the toss of a coin that can be biased
   $P(heads) = p, \quad P(tails) = 1 - p$, parameter is $p$.

2. <u>Binomial RV</u>—model number of heads, **k**, of $n$ biased coin
   tosses.

   $$P(x = k) = \binom{n}{k} p^k (1-p)^{n-k}$$

# e.g. distributions of discrete RVs

1. Bernouilli RV—model the toss of a coin that can be biased
   $P(heads) = p,\quad P(tails) = 1 - p$, parameter is $p$.

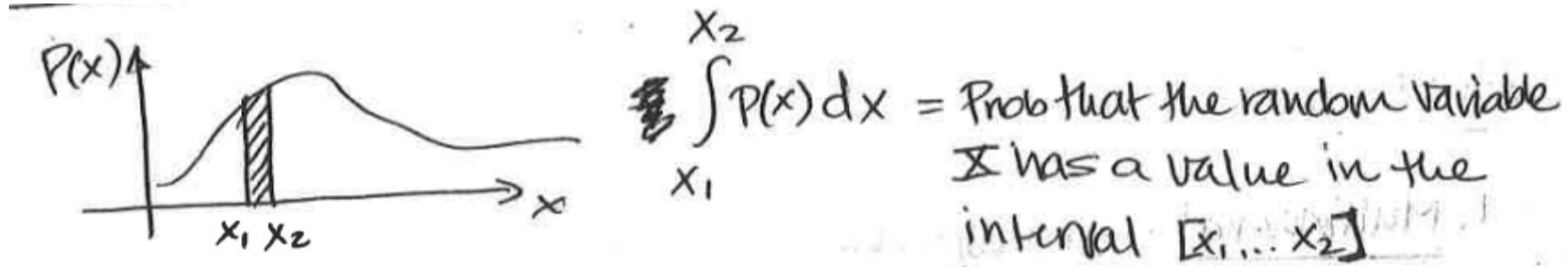2. Binomial RV—model number of heads, $\mathbf{k}$, of $n$ biased coin tosses.

$$P(x=k) = \binom{n}{k} p^k (1-p)^{n-k}$$

3. Poisson RV– model number of mutations, $k$, occurring in a cell population with mean mutation rate, $\lambda$, over fixed time interval
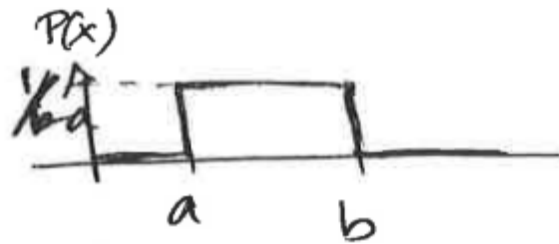
$$P(x=k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

# Distributions of continuous RVs

Continuous RVs have a Probability Density Function



$\int\limits_{X_1}^{X_2} P(x)\,dx$ = Prob that the random variable $X$ has a value in the interval $[X_1 ... X_2]$

Examples!

1. Uniform —

Parameters: $a, b$

2. Gaussian $\quad P(x) = \dfrac{1}{\sqrt{2\pi}\,\sigma} \exp\left(-\tfrac{1}{2}\dfrac{(x-\mu)^2}{\sigma^2}\right) \quad$ Parameters: $\mu, \sigma$

$$X \sim N(\mu, \sigma^2)$$

# Multivariate distributions

Space of outcomes is a vector instead of a scalar:
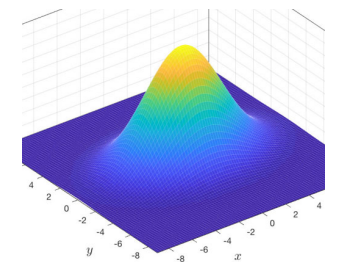
Multinomial (generalization from binomial):

- urn with balls of different colors.

- Pick a ball at random.

- $p_1$ it is green, $p_2$ it is blue and $p_3$ it is red

Multivariate Gaussian:

- Mean is a vector, and variance becomes covariance.

- Will learn more about this next lecture.

# The basic set-up of MLE

- Given data $D = \{x_i\}_{i=1}^{N}$ for $x_i \in R^d$
- Assume a set (family) of distributions on $R^d$, $\{p_\theta(x)|\theta \in \Theta\}$.

Same as
$p(x|\theta)$

e.g mean $(\mu)$ and variance $(\sigma^2)$ for $x \in \mathbb{R}^1$

# The basic set-up of MLE

- Given data $D = \{x_i\}_{i=1}^N$ for $x_i \in R^d$

- Assume a set (family) of distributions on $R^d$, $\{p_\theta(x) | \theta \in \Theta\}$.

- Assume $D$ contains samples from one of these distributions:

$$x_i \sim p_{\hat{\theta}}(x)$$

- This assumes that each element of $D$ is *identically and independently distributed* (iid).

*Same as $p(x|\theta)$*

# The basic set-up of MLE

- Given data $D = \{x_i\}_{i=1}^{N}$ for $x_i \in R^d$
- Assume a set (family) of distributions on $R^d$, $\{p_\theta(x) | \theta \in \Theta\}$.

  *e.g mean ($\mu$) and variance ($\sigma^2$) for $x \in R^1$*

- Assume $D$ contains samples from one of these distributions:

$$x_i \sim p_{\hat{\theta}}(x)$$

- This assumes that each element of $D$ is *identically and independently distributed* (iid).

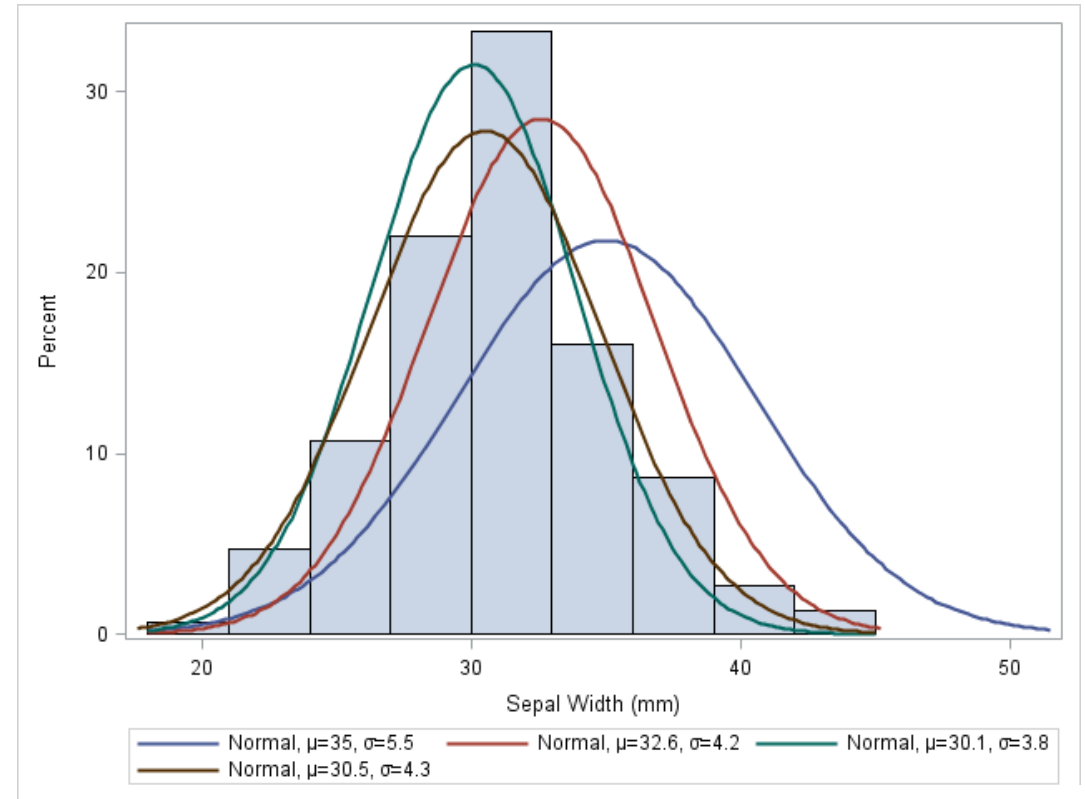Goal of MLE: "learn"/estimate the value of $\theta$ that defines the distribution from which the data came.

Definition: $\theta_{MLE}$ is a MLE for $\theta$ with respect to the data and set of distributions, if $\theta_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} \, p(D | \theta)$.

*"likelihood function"*

*function of $\Theta$.*

# The basic set-up of MLE

$$\theta_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} \, p(D|\theta)$$



Normal, μ=35, σ=5.5 — Normal, μ=32.6, σ=4.2 — Normal, μ=30.1, σ=3.8
Normal, μ=30.5, σ=4.3

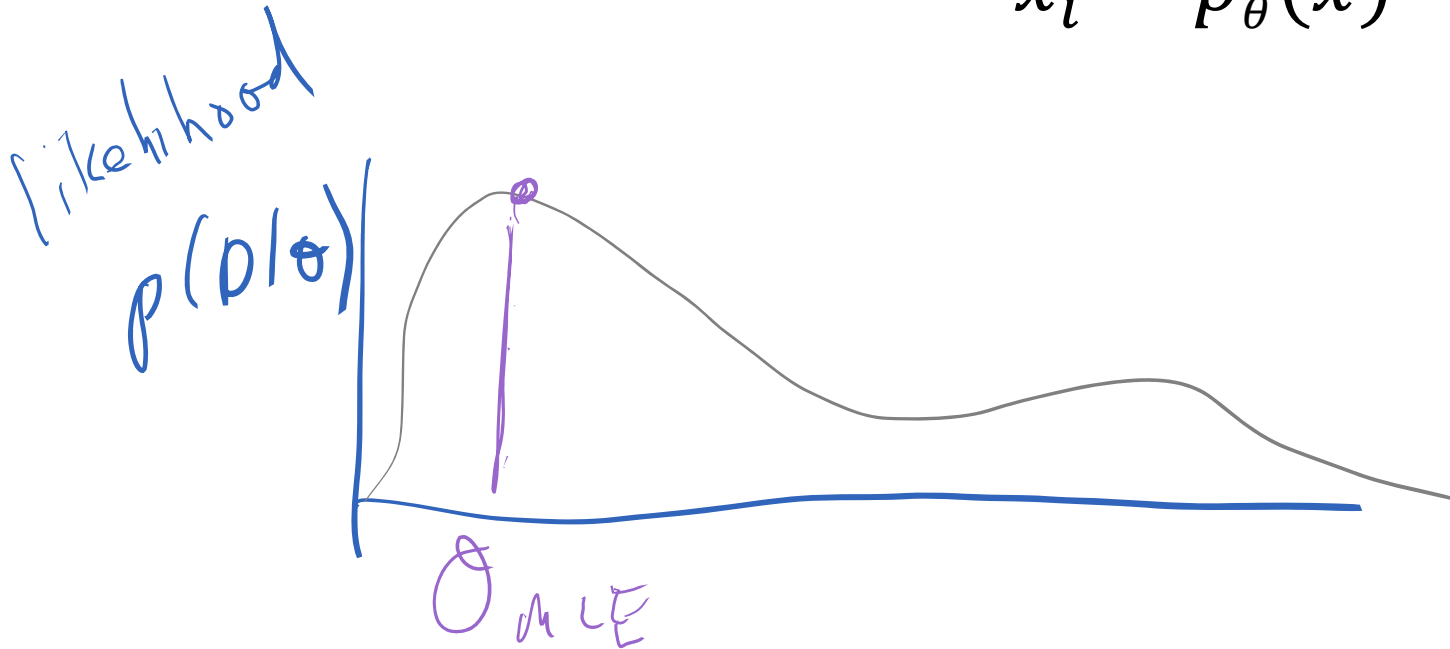$D = \{x_i\}_{i=1}^N = \{20.1, 33.8, 34.6, 36.2, \ldots\}$

Note that $p(D|\theta) = p(\{x_i\}_{i=1}^N|\theta) = \prod_{i=1}^N p(x_i|\theta)$

because iid

# The basic set-up of MLE

- Given data $D = \{x_i\}_{i=1}^{N}$ for $x_i \in R^d$

- Assume a set (family) of distributions on $R^d$, $\{p_\theta(x) | \theta \in \Theta\}$.

- Assume $D$ contains samples from one of these distributions:

$$x_i \sim p_{\hat{\theta}}(x)$$

$$\boxed{\theta_{MLE} = \underset{\theta \in \Theta}{\mathrm{argmax}} \, p(D|\theta)}$$

*Is there always one unique MLE parameter value?*

likelihood

$p(D|\theta)$

$\theta_{MLE}$

# Some properties of MLE

*eg* $N(x|\mu,\sigma)$
vs
$N(x|2+\mu,\sigma)$

- The MLE is a *consistent* estimator: meaning that as we get more and more data (drawn from one distribution in our family), then we converge to estimating the true value of $\boldsymbol{\theta}$ for $\boldsymbol{D}$.

- The MLE is *statistically efficient*: it's making good use of the data available to it ( "least variance" parameter estimates).

- The value of $p(\boldsymbol{D}|\boldsymbol{\theta_{MLE}})$ is invariant to re-parameterization.

- MLE can still yield a parameter estimate even when the data were not generated from that family (phew & caveat emptor).

# *e.g.* MLE for univariate Gaussian



- Arguments can be made from the Central Limit Theorem that height is normally distributed.
- Suppose you were given a set if height measurements, $\{x_i\}$, how would you derive the estimate for the mean and variance, using MLE?

# *e.g.* MLE for univariate Gaussian

Goal: $\theta_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} \, p(D|\theta)$ from set of data $D = \{x_i\}_{i=1}^{N}$

- Assume data are generated as $X \sim N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} exp - \frac{(x-\mu)^2}{2\sigma^2}$

- So assume MLE family of distributions, $p(X = x|\theta) = N(X|\mu, \sigma^2)$.

- Now our goal is to find $\theta_{MLE} = (\mu_{MLE}, \sigma^2_{MLE}) = \underset{\theta \in \Theta}{\operatorname{argmax}} \, p(D|\mu, \sigma^2)$.

- First step, write down the likelihood function:
  - $p(D|\theta) = p(x_1, x_2, \dots x_N|\mu, \sigma^2) = \prod_{i=1}^{N} p(x_i|\mu, \sigma^2)$.

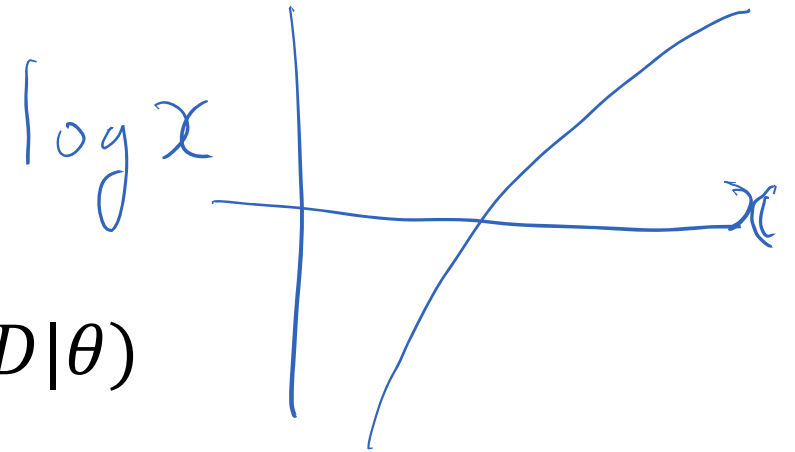- The product of the terms is a little inconvenient to work with.

# *e.g.* MLE for univariate Gaussian

- Likelihood: $p(x_1, x_2, \ldots x_N | \mu, \sigma^2) = \prod_{i=1}^{N} p(x_i | \mu, \sigma^2)$.

- The *log likelihood ("LL")* is a monotonically increasing function of the likelihood.

$$\log p(D|\theta) = \sum_{i=1}^{N} \log p(x_i | \mu, \sigma^2)$$

- Therefore $\theta_{MLE} = \underset{\theta \in \Theta}{\mathrm{argmax}}\, p(D|\theta) = \underset{\theta \in \Theta}{\mathrm{argmax}}\, \log p(D|\theta)$

# *e.g.* MLE for univariate Gaussian

- Now we have a concrete optimization problem to work with:

$$\mu_{MLE}, \sigma^2_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} \log p(D|\theta) = \underset{\mu,\sigma^2}{\operatorname{argmax}} \sum_{i=1}^{N} \log p(x_i|\mu,\sigma^2)$$

- How will we solve this optimization problem?
- Find a setting of the parameters for which the partial derivatives are 0 (*i.e.*, a stationary point).
- Then check whether the setting is a maximum (negative second derivative), a minimum, etc. (first year calculus).
- (if #params>1, check if Hessian is negative definite; for 1D Gaussian, Hessian is diagonal, so can check each separately).

# e.g. MLE for univariate Gaussian

- Find the setting of the parameters that set the partial derivatives to zero:

$$\mu_{MLE}, \sigma^2_{MLE} = \underset{\theta \in \Theta}{\text{argmax}} \log p(D|\theta) = \underset{\mu, \sigma^2}{\text{argmax}} \sum_{i=1}^{N} \log p(x_i|\mu, \sigma^2)$$

- Lets expand out so we can take the derivative:

$$\frac{\partial}{\partial \mu} \sum_{i=1}^{N} \log p(x_i|\mu, \sigma^2) = \sum_{i} \frac{\partial}{\partial \mu} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu)^2\right)\right]$$

# e.g. MLE for univariate Gaussian

- Find the setting of the parameters that set the partial derivatives to zero:

$$\mu_{MLE}, \sigma^2_{MLE} = \underset{\theta \in \Theta}{\mathrm{argmax}} \log p(D|\theta) = \underset{\mu, \sigma^2}{\mathrm{argmax}} \sum_{i=1}^{N} \log p(x_i|\mu, \sigma^2)$$

- Lets expand out so we can take the derivative:

$$\frac{\partial}{\partial \mu} \sum_{i=1}^{N} \log p(x_i | \mu, \sigma^2) = \sum_i \frac{\partial}{\partial \mu} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left( -\frac{1}{2\sigma^2}(x_i - \mu)^2 \right) \right]$$

$$= \sum_i \frac{\partial}{\partial \mu} \left[ -\log(\sqrt{2\pi\sigma^2}) - \frac{1}{2\sigma^2}(x_i - \mu)^2 \right]$$

$$= \sum_i \left[ 0 + \frac{1}{\sigma^2}(x_i - \mu) \right] \overset{\text{set to zero}}{\Longrightarrow} \sum_i x_i = \sum_i \mu$$

$$\sum_i x_i = N\mu$$

$$\Rightarrow \mu \quad \frac{\sum_i x_i}{N}$$

# *e.g.* MLE for univariate Gaussian

$\dfrac{\partial^2 (LL)}{\partial \mu^2} =$

- Find the setting of the parameters that set the partial derivatives to zero:

$$\mu_{MLE}, \sigma^2_{MLE} = \underset{\theta \in \Theta}{\mathrm{argmax}} \log p(D|\theta) = \underset{\mu, \sigma^2}{\mathrm{argmax}} \sum_{i=1}^{N} \log p(x_i | \mu, \sigma^2)$$

- Lets expand out so we can take the derivative:

$$\frac{\partial}{\partial \mu} \left( \sum_{i=1}^{N} \log p(x_i | \mu, \sigma^2) \right) = \sum_i \frac{\partial}{\partial \mu} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left( -\frac{1}{2\sigma^2}(x_i - \mu)^2 \right) \right]$$

$$= \sum_i \frac{\partial}{\partial \mu} \left[ -\log(\sqrt{2\pi\sigma^2}) - \frac{1}{2\sigma^2}(x_i - \mu)^2 \right]$$

$$= \sum_i \left[ 0 + \frac{1}{\sigma^2}(x_i - \mu) \right] \implies \text{set to zero}$$

$$\sum_i x_i = \sum_i \mu \implies \sum_i x_i = N\mu$$

$$\implies \mu \quad \frac{\sum_i x_i}{N}$$

# *e.g.* MLE for univariate Gaussian

$$\frac{d^2(LL)}{d\mu^2} = \sum_i \frac{1}{\sigma^2} \cdot (-1) = -\frac{N}{\sigma^2} < 0 \implies \text{min.}$$

$$\frac{d}{d\mu}\left(\sum_{i=1}^{N} \log p(x_i | \mu, \sigma^2)\right) = \sum_i \frac{d}{d\mu} \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu)^2\right)\right] \quad \text{MLE}$$

$$= \sum_i \frac{d}{d\mu}\left[-\log(\sqrt{2\pi\sigma^2}) - \frac{1}{2\sigma^2}(x_i - \mu)^2\right]$$

$$= \sum_i \left[0 + \frac{1}{\sigma^2}(x_i - \mu)\right] \xrightarrow{\text{set to zero}} \sum_i x_i = \sum_i \mu \implies \sum_i x_i = N\mu$$

$$\implies \mu \quad \frac{\sum_i x_i}{N}$$

# *e.g.* MLE for univariate Gaussian

$$N(x|\mu,\sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} exp - \frac{(x-\mu)^2}{2\sigma^2}$$

$$\mu_{MLE}, \sigma^2_{MLE} = argmax \sum_{i=1}^{N} \log N(x_i|\mu,\sigma^2)$$

- Again, but this time for $\sigma^2$:

$$\frac{\partial \mathcal{LL}}{\partial \sigma^2} = \frac{\partial}{\partial \sigma^2}\left(-\frac{N}{2}\log(2\pi\sigma^2) - \frac{1}{2\sigma^2}\sum_{n=1}^{N}(x_n-\mu)^2\right)$$

$$= -\frac{N}{2}\cdot\frac{\partial}{\partial\sigma^2}\left(\log(2\pi\sigma^2)\right) + \frac{\partial}{\partial\sigma^2}\left(-\frac{1}{2\sigma^2}\sum_{n=1}^{N}(x_n-\mu)^2\right)$$

$$= -\frac{N}{2}\cdot\frac{1}{2\pi\sigma^2}\cdot 2\pi + \frac{\partial}{\partial\sigma^2}\left(-\frac{1}{2\sigma^2}\sum_{n=1}^{N}(x_n-\mu)^2\right)$$

$$= -\frac{N}{2\sigma^2} + \sum_{n=1}^{N}\left(-\frac{1}{2}\cdot -1\cdot(\sigma^2)^{-2}\cdot 1\cdot(x_n-\mu)^2\right)$$

$$= -\frac{N}{2\sigma^2} + \sum_{n=1}^{N}\left(\frac{1}{2\sigma^4}\cdot(x_n-\mu)^2\right)$$

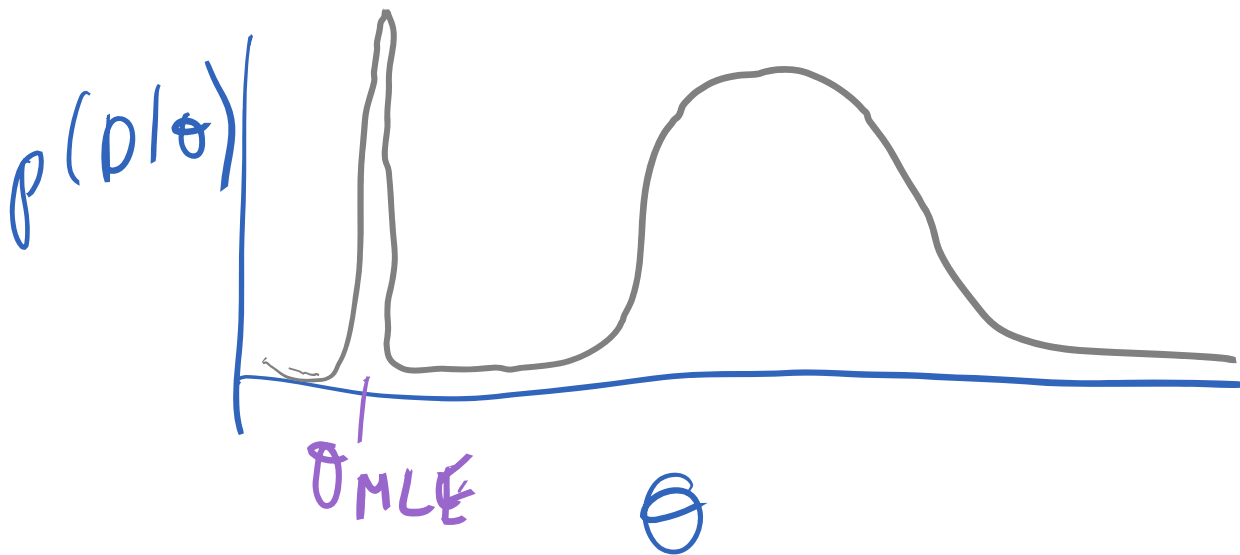$$0 = \frac{1}{2\sigma^2}\left(-N + \frac{1}{\sigma^2}\sum_{n=1}^{N}(x_n-\mu)^2\right)$$

$$0 = -N + \frac{1}{\sigma^2}\sum_{n=1}^{N}(x_n-\mu)^2$$

$$\sigma^2 = \frac{1}{N}\sum^{N}(x_n-\mu)^2$$
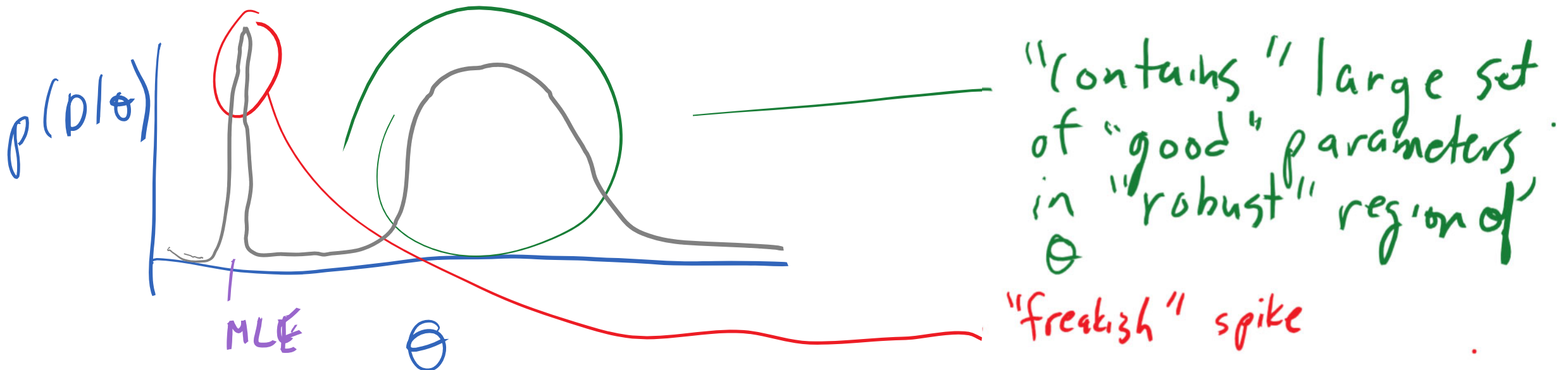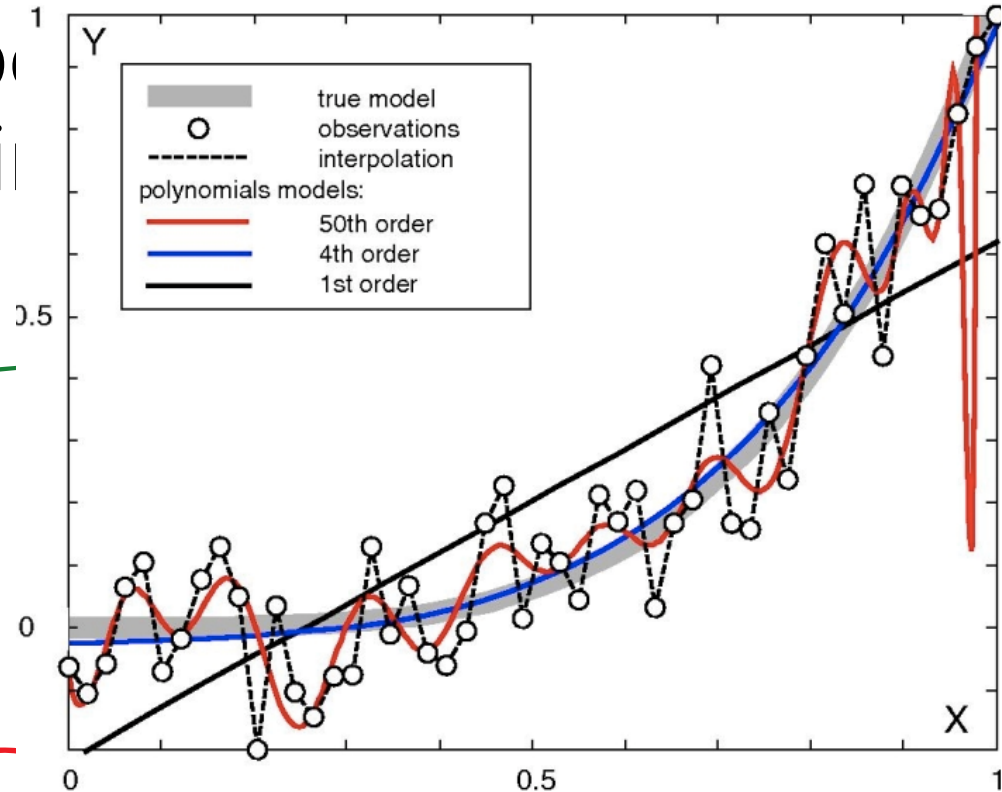
$$\sigma^2_{MLE}$$

# MLE yields a "point estimate" of our parameter

- When we perform MLE, we get just one single estimate of the parameter, $\theta$, rather than a distribution over it which captures uncertainty.

- In Bayesian statistics, we obtain a (posterior) distribution over $\theta$. We will touch more on this in a few lectures.
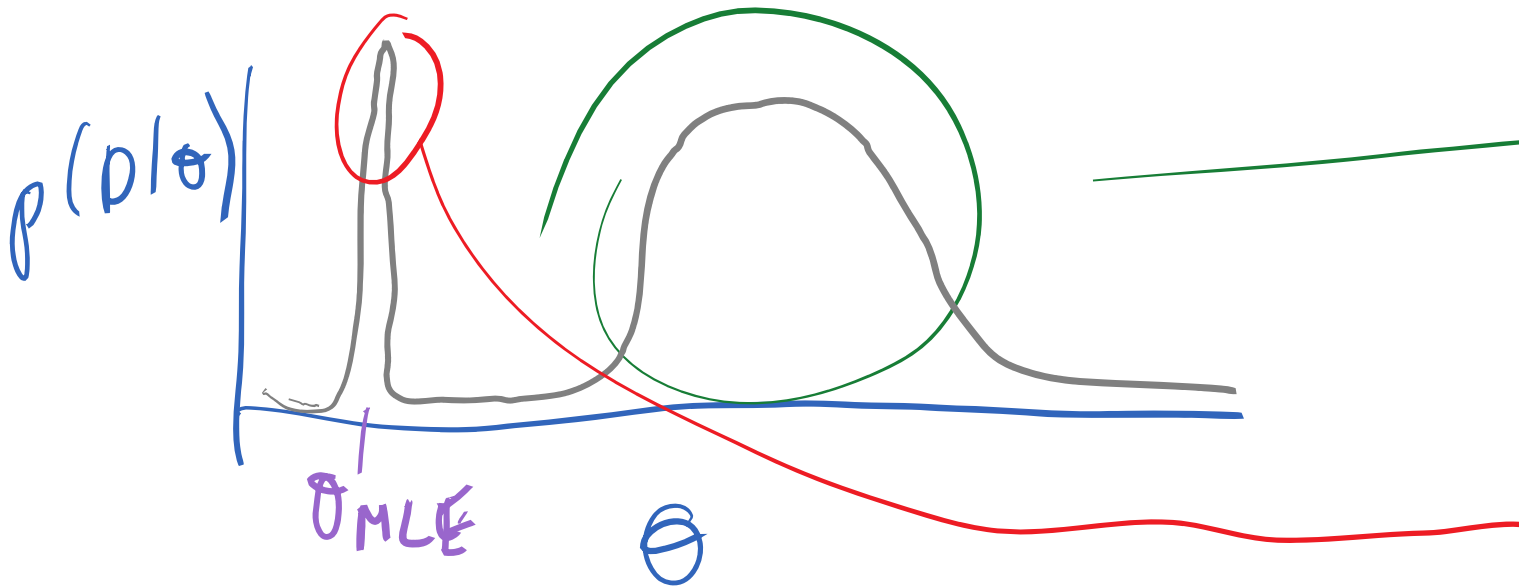
$p(D|\theta)$

$\theta_{MLE}$

$\theta$

# MLE yields a "point estimate" of our parameter

- When we perform MLE, we get just <u>one single estimate of the parameter, $\theta$,</u> rather than a distribution over it which captures uncertainty.

- In Bayesian statistics, we obtain a (posterior) distribution over $\theta$. We will touch more on this in a few lectures.



$p(D|\theta)$

MLE

$\theta$

"contains" large set of "good" parameters in "robust" region of $\theta$
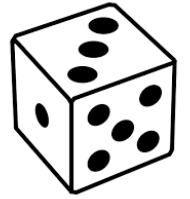
"freakish" spike

# MLE yields a "point estimate" of our parameter

- When we perform MLE, we get just one single estimate of the parameter, $\theta$, rather than a distribution over it which captures uncertainty.

- In Bayesian statistics, we obtain a (p[ ] over $\theta$. We will touch more on this i[ ]

# e.g. MLE for the multinomial distribution

- Consider a six-sided die that we will roll, and we want to know the probability of each side of the die turning up ($\theta = \theta_1 \dots \theta_6$).
- Assume we have observed $N$ rolls, with RV, $X \sim p_\theta(X)$.
- We write that $P(X = k|\theta) = \theta_k$ (when $k^{th}$ side faced up).
- Lets use MLE to estimate these parameters.
- First, since one side must always face up, we know that $1 = \sum_k \theta_k$.
- Second, we can write $P(X = x|\theta) = \theta_x$ (pick off the right parameter).
- Now we write the likelihood:

$$n_k \equiv \left| \{ i \,|\, x_i = k \} \right|$$

$$P(D|\theta) = p(x_1, \dots x_N|\theta) = \prod_{i=1}^{N} p(x_i|\theta) = \prod_{i=1}^{N} \prod_{k=1}^{6} \theta_k^{I[x_i=k]} = \prod_{k=1}^{6} \theta_k^{\sum_i^N I[x_i=k]} = \prod_{k=1}^{6} \theta_k^{n_k}$$
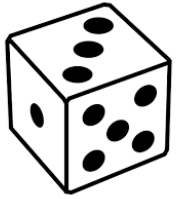
Now our MLE problem becomes:

constrained optimization

$$\theta_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} \log p(D|\theta) = \underset{\theta \in \{\Theta | 1 = \sum_k \theta_k\}}{\operatorname{argmax}} \sum_{k=1}^{6} \log \theta_k^{n_k}$$

# *e.g.* MLE for the multinomial distribution

Have a constrained optimization problem:

$$\theta_{MLE} = \underset{\theta \in \Theta}{\mathrm{argmax}} \log p(D|\theta) = \underset{\theta \in \{\Theta | 1 = \sum_k \theta_k\}}{\mathrm{argmax}} \sum_{k=1}^{6} \log \theta_k^{n_k}$$
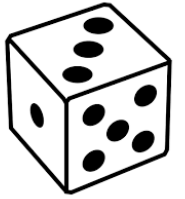
*constrained optimization*

What is one technique you should have learned in first year calculus to solve this?

The technique of Lagrange multipliers:

$$J(\theta, \lambda) = \log p(D|\theta) + \lambda(1 - \sum_k \theta_k)$$ (look for stationary points wrt $\theta, \lambda$)

# *e.g.* MLE for the multinomial distribution

$$J(\theta, \lambda) = \log p(D|\theta) + \lambda(1 - \sum_k \theta_k) = \sum_{k=1}^{6} \log \theta_k^{n_k} + \lambda(1 - \sum_k \theta_k)$$

1. $\frac{\partial J}{\partial \lambda} = 0 \Rightarrow 1 = \sum_k \theta_k$ (we just get the constraint back)
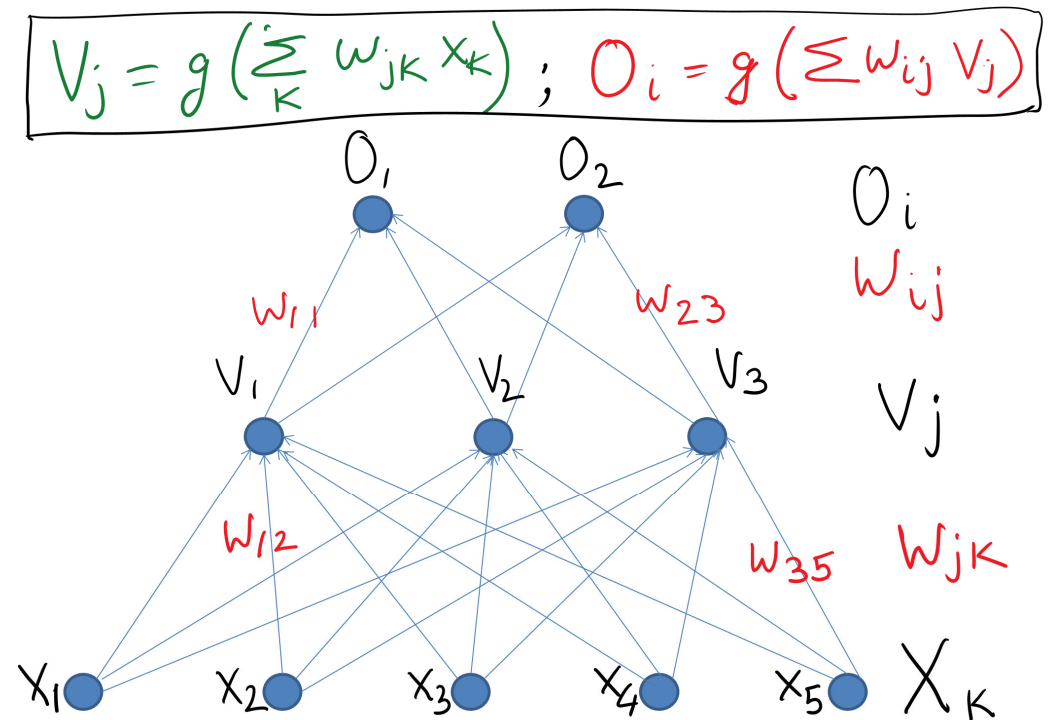
2. $\frac{\partial J}{\partial \theta_k} = \frac{\partial}{\partial \theta_k} \sum_{k=1}^{6} \log \theta_k^{n_k} - \frac{\partial}{\partial \theta_k} \lambda \theta_k = \frac{n_k}{\theta_k} - \lambda = 0 \Rightarrow \theta_k = \frac{n_k}{\lambda}.$

3. Lets plug this into 1), $1 = \sum_k \theta_k = \sum_k \frac{n_k}{\lambda} \Rightarrow \lambda = \sum_k n_k = N.$

4. All together then, $\theta_k = \frac{n_k}{N}.$

# Doing MLE requires optimization

$$\theta_{MLE} = \operatorname*{argmax}_{\theta \in \Theta} \log p(D|\theta)$$

- For Gaussian, multinomial (and more), the MLE can be obtained in closed form by setting the derivative to zero.
- What if we had a model such as Prof. Malik mentioned in the first lecture?

- Here, we need *iterative optimization* (can take entire classes on special cases of this (e.g. Convex Optimization). More later.

# Prof. Malik in first lecture:

- Mentioned that a good loss to estimate parameters is the *cross-entropy* (rather than the likelihood).
- So why are we teaching you MLE?! They are equivalent.

Training a single layer neural network
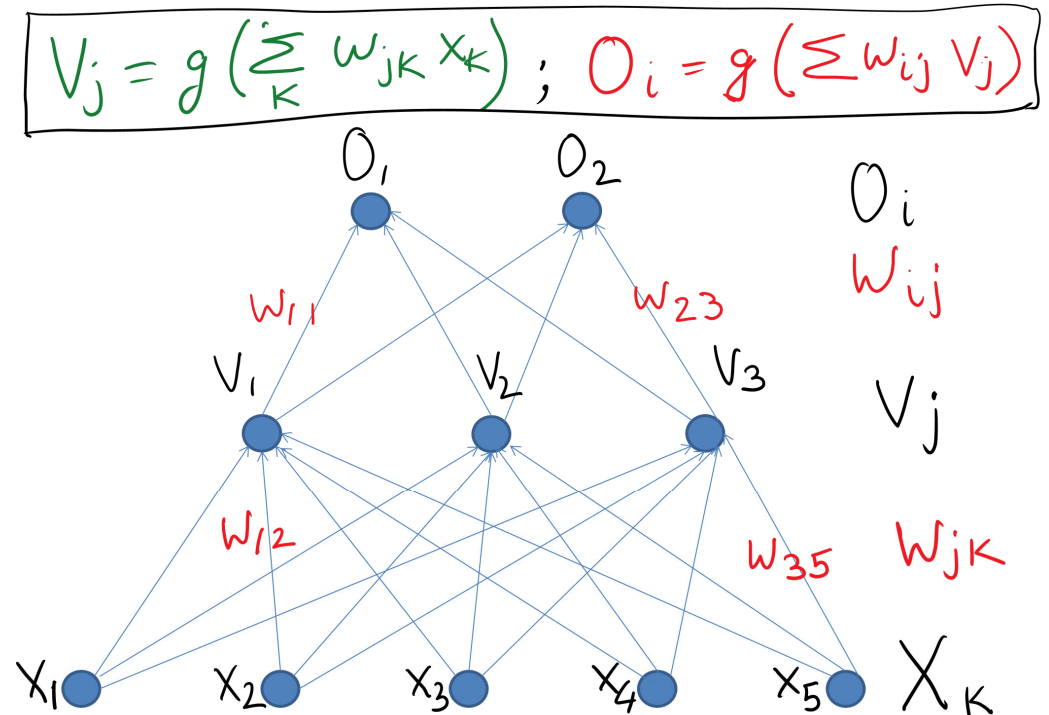
- A good choice of loss function is the cross entropy

$$L = - \sum_{\text{input data}} \left( y_i \ln O_i + (1-y_i) \ln(1-O_i) \right)$$

- We model the activation function *g* as a sigmoid

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

- Finding w reduces to logistic regression!

We can use STOCHASTIC GRADIENT DESCENT.

$$V_j = g\left(\sum_K w_{jK} x_K\right) ; \quad O_i = g\left(\sum w_{ij} V_j\right)$$

# Relationship between likelihood, cross-entropy, *etc.*

- The *cross-entropy* is a term from *information* theory.
- To understand the connection between MLE and maximizing the cross-entropy, we need to know some concepts from information theory:
  1. Entropy
  2. Cross-entropy
  3. KL-divergence (relative entropy).

# Entropy: a measure of expected surprise

Think about a flipping a coin once, and how surprised you would be at observing a head.

$p(head) = 0.5$

$p(head) = 0$

$p(head) = 1$

$p(head) = 0.01$

# Entropy: a measure of expected surprise

- The "surprise" of observing that a discrete random variable $Y$ takes on value $k$ is:

$$log \frac{1}{P(Y = k)} = -\log(P(Y = k))$$

- As $P(Y = k) \to 0$, the surprise of observing $k$ approaches $\infty$.
- As $P(Y = k) \to 1$, the surprise of observing $k$ approaches $0$.
- The entropy of the distribution of $Y$ is the *expected surprise*:

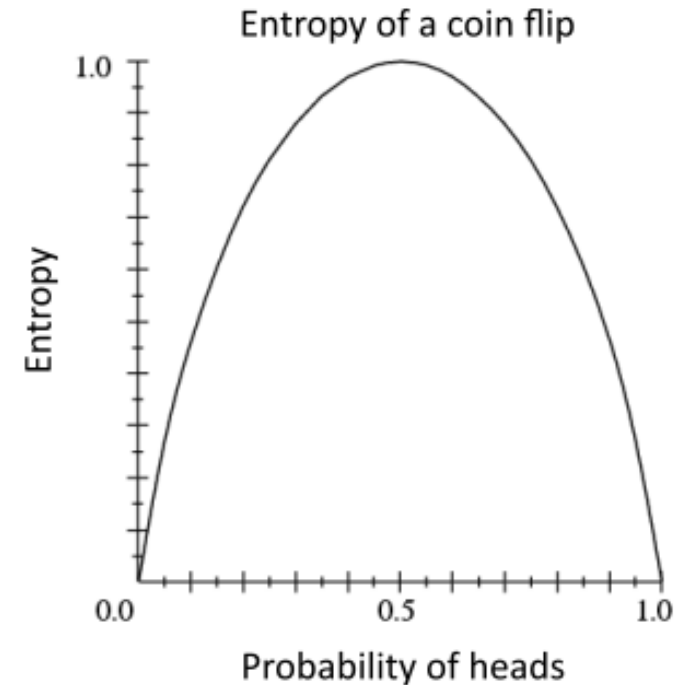$$H(Y) \equiv E_Y[-\log P(Y = k)] = \sum_k P(Y = k) \log P(Y = k)$$

# Entropy example: flipping a coin

$$H(Y) = -\sum_{i=1}^{\kappa} P(Y = y_i) \log_2 P(Y = y_i)$$

P(Y=t) = 5/6

P(Y=f) = 1/6

H(Y) = - 5/6 log$_2$ 5/6 - 1/6 log$_2$ 1/6

= 0.65

Entropy of a coin flip



Entropy (y-axis) vs Probability of heads (x-axis)
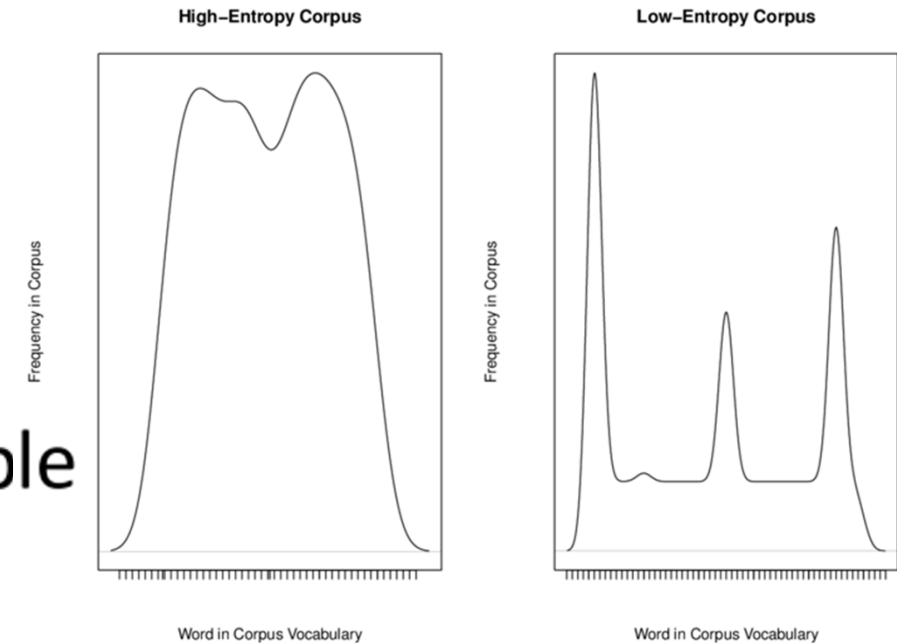
# Entropy of a random variable $Y$:

"High Entropy"

– Y is from a uniform like distribution

– Flat histogram

– Values sampled from it are less predictable

"Low Entropy"

– Y is from a varied (peaks and valleys) distribution

– Histogram has many lows and highs

– Values sampled from it are more predictable



High–Entropy Corpus

Low–Entropy Corpus

Frequency in Corpus

Word in Corpus Vocabulary

*https://www.researchgate.net/figure/Hypothetical-distributions-of-term-frequency-in-high-and-low-entropy-corpora_fig1_305417514*

(Slide from Vibhav Gogate)

# From Entropy to *Relative Entropy*

- Also called the Kullback-Leibler (KL) Divergence.
- Measures how much one distribution diverges from another.
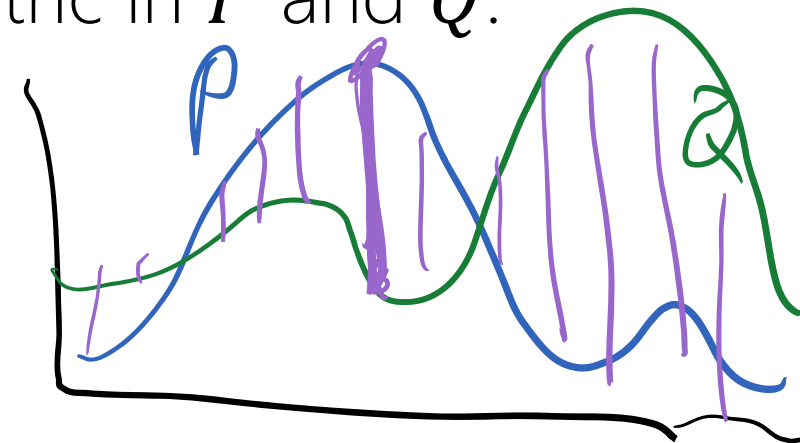- For discrete probability distributions, $P$ and $Q$, it is defined as:

$$D_{KL}(P||Q) = \sum_x P(x) \ln \frac{P(x)}{Q(x)}$$

- Not a true distance metric because not symmetric in $P$ and $Q$:

$$D_{KL}(P||Q) \neq D_{KL}(Q||P)$$

## Properties of KL Divergence
- $\mathrm{KL}(p||q) \geq 0$
- $\mathrm{KL}(p||q) = 0$ if and only if $p = q$

# From Relative Entropy to Cross-Entropy (then to MLE!)

$$D_{KL}(P||Q) = \sum_x P(x) log \frac{P(x)}{Q(x)}$$

$$= E_{P(x)}\left[log \frac{1}{Q(X)}\right] - E_{P(x)}\left[log \frac{1}{P(X)}\right]$$

# From Relative Entropy to Cross-Entropy (then to MLE!)

$$D_{KL}(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

$$= E_{P(x)}\left[\log \frac{1}{Q(X)}\right] - E_{P(x)}\left[\log \frac{1}{P(X)}\right]$$

$$= H(P,Q) - H(P)$$

cross-entropy          entropy

- Consider data, $D$ where $x_i \sim \hat{p}_{data}$ and a model with params $\theta$, $p(x|\theta)$.
- If minimizing the KL divergence (instead of MLE),
$$argmin_\theta D_{KL}(\hat{p}_{data}||p(x|\theta)) = )) =$$

# From Relative Entropy to Cross-Entropy (then to MLE!)

$$D_{KL}(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}$$

$$= E_{P(x)}\left[\log \frac{1}{Q(X)}\right] - E_{P(x)}\left[\log \frac{1}{P(X)}\right]$$

$$= H(P, Q) - H(P)$$

cross-entropy

entropy

no dependence on θ

- Consider data, $D$ where $x_i \sim \hat{p}_{data}$ and a model with params $\theta$, $p(x|\theta)$.
- If minimizing the KL divergence (instead of MLE),

$$argmin_\theta D_{KL}(\hat{p}_{data}||p(x|\theta)) = )) = argmin_\theta H(\hat{p}_{data}, p(x|\theta)) + H(\hat{p}_{data})$$

$$= argmax E_{\hat{p}_{data}}[\log p(x|\theta)]$$

negative cross-entropy

# From Relative Entropy to Cross-Entropy (then to MLE!)

$$D_{KL}(P||Q) = \sum_x P(x) log \frac{P(x)}{Q(x)}$$

$$= E_{P(x)}\left[log \frac{1}{Q(X)}\right] - E_{P(x)}\left[log \frac{1}{P(X)}\right]$$

$$= H(P,Q) - H(P)$$

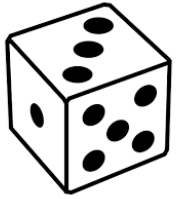cross-entropy      entropy

MLE problem

- Consider data, $D$ where $x_i \sim \hat{p}_{data}$ and a model with params $\theta$, $p(x|\theta)$.
- If minimizing the KL divergence (instead of MLE),

$$argmin_\theta D_{KL}(\hat{p}_{data}||p(x|\theta)) = )) = argmin_\theta H(\hat{p}_{data}, p(x|\theta)) + H(\hat{p}_{data})$$

$$= argmax E_{\hat{p}_{data}}[\log p(x|\theta)] = argmax \sum_i^N \log p(x_i|\theta).$$

# From Relative Entropy to Cross-Entropy (then to MLE!)

- Performing MLE maximizes the likelihood function.
- This is equivalent to maximizing the cross-entropy.
- And equivalent to minimizing the KL-divergence (aka relative entropy).

# Extra

# *e.g.* MLE for the multinomial distribution

$J(\theta, \lambda) = \log p(D|\theta) + \lambda(1 - \sum_k \theta_k)$ (look for stationary points wrt $\theta, \lambda$)

1. $\frac{\partial J}{\partial \lambda} = 0 \rightarrow 1 = \sum_k \theta_k$ (we just get the constraint back)

2. $\frac{\partial J}{\partial \theta_k} = \frac{\partial}{\partial \theta_k} \sum_{k=1}^{6} \log \theta_k^{n_k} - \frac{\partial}{\partial \theta_k} \lambda \theta_k = \frac{n_k}{\theta_k} - \lambda = 0 \rightarrow \theta_k = \frac{n_k}{\lambda}$.

3. Lets plug this into 1: $\sum_k \theta_k = 1 = \sum_k \frac{n_k}{\lambda} \rightarrow \lambda = \sum_k n_k = N$.

4. All together then, $\theta_k = \frac{n_k}{N}$.

This is a stationary point. But is it a maximum? Could check Hessian, but lets instead consider our know equivalence

$$D_{KL}(p_{data}||p(x|\theta)) = \sum_{k=1}^{6} P_{data}(X = k) \log \frac{P_{data}(X=k)}{P(X=k|\theta)}$$

*Handwritten annotations:*
$\log(1) = 0$
$\Rightarrow D_{KL} = 0$
but $D_{KL} \geq 0$
So we have minimize the $D_{KL}$
$n_k/N$
maximized log likelihood !
$\Rightarrow \theta_k = \frac{n_k}{N}$