

Machine Learning Notes



Probability:

1. marginal probability = $\int (\text{probability A} , \text{probability B}) * d(\text{probability B}) = \int (\text{probability A} | \text{probability B}) * \text{probability B} * d(B)$
2. $\sum P(A_i | B) = 1$
3. $p(Y | X) = \int p(Y | Z, X) p(Z | X) dZ$
4. posterior probability = prior probability * marginal of the prior probability / marginal of the posterior probability
5. $P(AB|C) P(C) = P(A|BC) P(BC)$
thus, $P(AB|C) = P(A|BC) P(BC) / P(C) = P(A|BC) P(B|C)P(C) / P(C) = P(A|BC) P(B|C)$

<https://www.youtube.com/watch?v=jY2E6ExLxaw>

7. Taking the expectation over D, based on the expectation linearity,
 $E[y(x;D) - E[y(x;D)]] = E[y(x;D)] - E[E[y(x;D)]] = E[y(x)] - E[E[y(x)] * 1] = E[y(x)] - E[y(x)] * E[1] = 0$
8. Probability model is always working on random variables
9. For random variables X and Y, X is independent of Y if and only if:
 - $P(X, Y) = P(X) * P(Y)$
 - $P(X | Y) = P(X)$
 - $P(Y | X) = P(Y)$
10. For random variables X, Y and Z, X is conditionally independent of Y given Z if and only if:
 - $P(X, Y | Z) = P(X | Z) * P(Y | Z)$
 - $P(X | Y, Z) = P(X | Z)$
 - $P(Y | X, Z) = P(Y | Z)$
 - $P(X, Y, Z) \propto \Phi_1(X, Z) * \Phi_2(Y, Z)$
11. Probability Mass Function (PMF for discrete random variable) =

Probability Density Function (PDF for concrete random variable) = Probability Distribution

12. If X_1, X_2, \dots, X_n are independent Bernoulli(p) random variables, then the random variable X defined by $X = X_1 + X_2 + X_3 + \dots + X_n$ has a Binomial(n, p) distribution

The law of large numbers:

Let $X_1, X_2, X_3 \dots$ be a set of independent and identically distributed random variables. It can be shown that:

$$\lim_{n \rightarrow \infty} P(|S_n/n - E(X)| > \epsilon) = 0$$

where $S_n = X_1 + X_2 + X_3 + \dots + X_n$

Central Limit Theorem:

Classical Definition:

Taking the mean of n samples that are i.i.d. sampled from the source distribution as a random variable S_n as the expectation of it according to **Law of Large Number**,

$$S_n := \frac{X_1 + \dots + X_n}{n}$$

As n gets larger, the random variable

$$\sqrt{n}(S_n - \mu)$$

approximates the normal distribution with mean 0 and variance σ^2 .

For large enough n , the distribution of S_n is close to the normal distribution with mean μ and variance σ^2 / n .

Lindeberg–Lévy Definition:

Let X_1, X_2, \dots, X_n be i.i.d random variables with expected value $E[X_i] = \mu < \infty$ and variance $0 < \text{Var}(X_i) = \sigma^2 < \infty$. Then, the random variable

$$\sqrt{n} \left(\left(\frac{1}{n} \sum_{i=1}^n X_i \right) - \mu \right) \xrightarrow{d} N(0, \sigma^2).$$

converges in distribution to the standard normal random variable as n goes to infinity, that is

$$\lim_{n \rightarrow \infty} \Pr[\sqrt{n}(S_n - \mu) \leq z] = \Phi\left(\frac{z}{\sigma}\right),$$

where $\Phi(x)$ is the standard normal distribution.

As i goes infinity,

$$\text{Var}(\mu) = \text{Var}((1/n) \sum_i X_i) = (1/n^2) \sum_i \text{Var}(X_i) = \sigma^2 / n$$

The variance of the mean decreases when i increases.

Why is central limit theorem important ?

Many statistics have distributions that are approximately normal for large sample sizes, even when we are sampling from a distribution that is not normal.

<https://www.khanacademy.org/math/ap-statistics/sampling-distribution-ap/sampling-distribution-mean/v/sampling-distribution-of-the-sample-mean>
https://www.youtube.com/watch?v=Pujol1yC1_A

Law of the unconscious statistician (LOTUS) for continuous random variables:

if $g(X)$ is a function of random variable X , then $g(X)$ is also a random variable, thus,

$$E[g(X)] = \int g(X) f(X) dX,$$

where f is the probability density function of X

Expectation:

1. Mean vs. Sample Mean

Let X_1, X_2, \dots, X_n be n i.i.d. random variables drawn from a distribution with mean μ and variance σ^2 , then

$$E[X_i] = \mu$$

While we define the random variable S_n is

$$S_n = (X_1 + X_2 + X_3 + \dots + X_n) / n = (1/n) \sum X_i$$

So the sample mean is

$$E[S_n] = E[(1/n) \sum X_i] = (1/n) \sum E[X_i] = (1/n) * n * \mu$$

$= \mu$

2. Let X and Y be random variables with finite expectations, $g(X)$ and $h(X)$ are also random variables, then

If $g(X) \geq h(X)$ for all real number x , then $E[g(X)] \geq E[h(X)]$

$E[aX + bY + c] = aE[X] + bE[Y] + c$ for any real number a, b, c

$E[X_1 + X_2 + \dots + X_n] = E[X_1] + E[X_2] + \dots + E[X_n]$, for any set of random variables X_1, X_2, \dots, X_n

3. Let X and Y be random variables with finite expectations, given X is known (which implies $g(X)$ is known as well), taken out what is known,

$E[g(X) h(Y) | X] = g(X) E[h(Y) | X]$

4. The law of iterated (total) expectation,

$E[Y] = \sum P(X = x_i) E[Y | X = x_i],$

seeing $X' = E[Y | X = x_i]$, with LOTUS, we have

$E[X'] = E[E[Y | X]] = \sum P(X = x_i) E[Y | X = x_i] = E[Y],$

thus,

$E[Y] = E[E[Y | X]]$

for more details,

<https://imai.princeton.edu/teaching/files/Expectation.pdf>

<https://www.probabilitycourse.com>

Variance / Covariance:

1. variance definition

variance is the expectation of the squared deviation of a random variable from its mean. Informally, it measures how far a set of (random) numbers are spread out from their average value.

$$\begin{aligned}\text{Var}(X) &= \sigma^2 \\ &= E[(x - E(x))^2] \\ &= E[X^2] - E[X]^2 \\ &= \int x^2 f(x) dx - \mu^2 \\ &= \text{Cov}(X, X)\end{aligned}$$

$$\text{Var}(aX) = a^2 \text{Var}(X)$$

like entropy, variance could sort of present the uncertainty (information) of a distribution.

<https://www.quora.com/Which-is-the-best-measure-of-uncertainty-variance-or-entropy-or-are-they-both-equivalent>

the entropy of a Gaussian distribution is proportional to the determinant of it.

2. Variance vs. Sample Variance

Let X_1, X_2, \dots, X_n be n i.i.d. random variables drawn from a distribution with mean μ and variance σ^2 , then

$$\text{Var}[X_i] = \sigma^2$$

While we define the random variable S_n is

$$S_n = (X_1 + X_2 + X_3 + \dots + X_n) / n = (1/n) \sum X_i$$

So the sample variance is

$$\text{Var}[S_n] = \text{Var}[(1/n) \sum X_i] = (1/n^2) \sum \text{Var}[X_i] = (1/n^2) * n * \sigma^2 = \sigma^2 / n$$

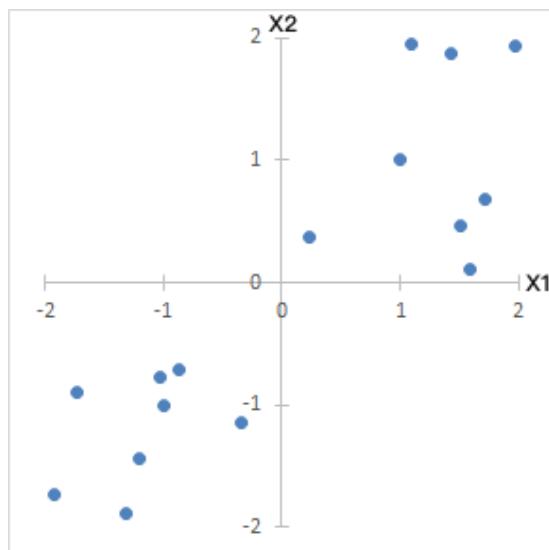
3. covariance definition

$$\text{Cor}(X, Y) = E[(x - \mu_x)(y - \mu_y)] = E[XY] - E[X]E[Y]$$

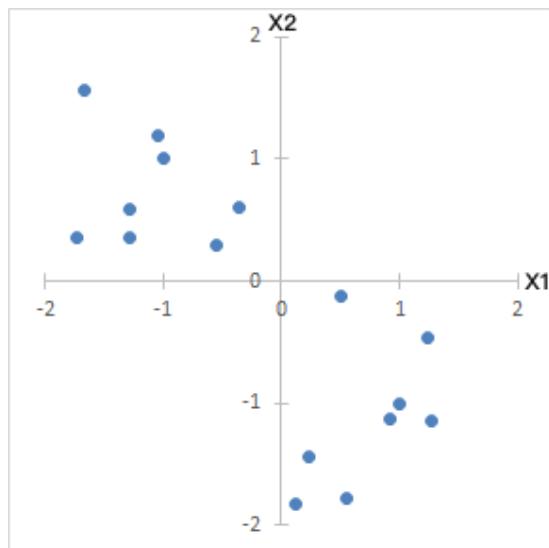
the covariance between X and Y indicates how the values of X and Y are more relative to each other.

an non-diagonal correlation matrix implies that the eigenvectors of the data distribution **are not aligned with the axes**.

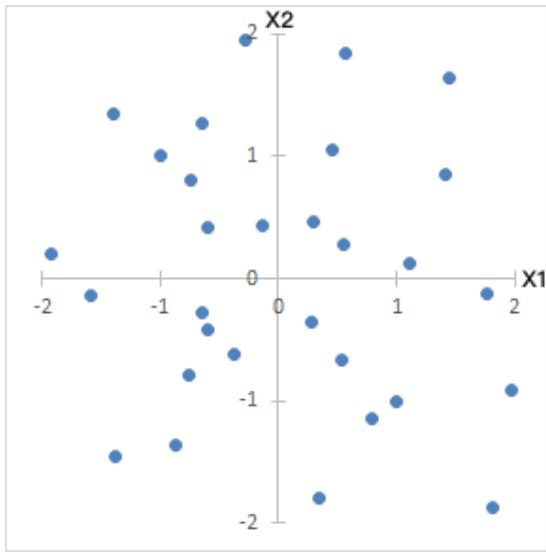
For example, Consider a 15×2 design matrix, which means 15 pairs of (x_1, x_2) data, with $E[x] = 0$,



$$Cov \approx \begin{bmatrix} \sum x_1^2 & \sum x_1 x_2 \\ \sum x_1 x_2 & \sum x_2^2 \end{bmatrix} \text{ where } \sum x_1 x_2 > 0$$



$$Cov \approx \begin{bmatrix} \sum x_1^2 & \sum x_1 x_2 \\ \sum x_1 x_2 & \sum x_2^2 \end{bmatrix} \text{ where } \sum x_1 x_2 < 0$$



$$Cov \approx \left[\begin{array}{cc} \sum x_1^2 & \sum x_1 x_2 \\ \sum x_1 x_2 & \sum x_2^2 \end{array} \right] \text{ where } \sum x_1 x_2 = 0$$

4. correlation coefficient

$$\rho(X, Y) = \frac{\text{Cov}(X - EX, Y - EY)}{\sigma X \sigma Y}, \quad -1 \leq \rho(X, Y) \leq 1$$

5. variance decomposition(law of total variance)

$$\text{Var}(X) = E[\text{Var}(X|Y)] + \text{Var}(E[X|Y]).$$

proof: Let $V = \text{Var}(X|Y)$, $Z = E[X|Y]$

$$V = E[X^2|Y] - (E[X|Y])^2 = E[X^2|Y] - Z^2$$

thus,

$$\begin{aligned} E[V] &= E[E[X^2|Y] - Z^2] \\ &= E[E[X^2|Y]] - E[Z^2] \quad \# \text{ expectation linearity} \\ &= E[X^2] - E[Z^2] \quad \# \text{ law of iterated expectation} \end{aligned}$$

next, we have

$$\text{Var}(Z) = E[Z^2] - (E[Z])^2$$

$$\begin{aligned}
&= E[Z^2] - (E[E[X|Y]])^2 \\
&= E[Z^2] - (E[X])^2 \text{ # law of iterated expectation}
\end{aligned}$$

Combining 2 equations above, we have

$$E[V] + \text{Var}(Z) = E[X^2] - (E[X])^2 = \text{Var}(X)$$

that is,

$$\text{Var}(X) = E[\text{Var}(X|Y)] + \text{Var}(E[X|Y])$$

the variance decomposition implies that

$$\text{Var}(X) \geq E[\text{Var}(X|Y)]$$

this states that when we condition on Y, the uncertainty (variance) of X reduces on average.

6. product of independent variables

$$\begin{aligned}
\text{Var}(XY) &= E[X^2Y^2] - E[XY]^2 \\
&= E[X^2]E[Y^2] - E[X]^2E[Y]^2 \\
&= E[X]^2E[Y^2] + E[Y]^2E[X^2] - 2E[X]^2E[Y]^2 + \\
&E[X^2]E[Y^2] - E[X]^2E[Y]^2 - E[X]^2E[Y^2] + E[X]^2E[Y]^2 \\
&= E[X]^2(E[Y^2] - E[Y]^2) + E[Y]^2(E[X^2] - E[X]^2) \\
&+ (E[X^2] - E[X]^2)(E[Y^2] - E[Y]^2) \\
&= E[X]^2\text{Var}(Y) + E[Y]^2\text{Var}(X) + \text{Var}(X)\text{Var}(Y)
\end{aligned}$$

7. if X_1, X_2, \dots, X_n are pairwise uncorrelated, i.e., $\rho(X_i, X_j) = 0$ when $i \neq j$, then

$$\text{Var}(X_1 + X_2 + \dots + X_n) = \text{Var}(X_1) + \text{Var}(X_2) + \dots + \text{Var}(X_n)$$

In case of 2 random variables X_1 and X_2 ,

$$\begin{aligned}
\text{Var}(X_1 + X_2) &= E[(X_1 + X_2)^2] - E[(X_1 + X_2)]^2 \\
&= E[X_1^2 + 2X_1X_2 + X_2^2] - (E[X_1] + E[X_2])^2 \\
&= E[X_1^2] + 2E[X_1X_2] + E[X_2^2] - E[X_1]^2 - \\
&\quad 2E[X_1]E[X_2] + E[X_2]^2 \\
&= (E[X_1^2] - E[X_1]^2) + (E[X_2^2] - E[X_2]^2) \\
&= \text{Var}(X_1) + \text{Var}(X_2)
\end{aligned}$$

8. variance of weighted sum random variables

$$\begin{aligned}
\text{Var}(c^T X) &= \text{Var}(\sum_i c_i X_i) \\
&= E[(\sum_i c_i X_i)^2] - E[(\sum_i c_i X_i)]^2
\end{aligned}$$

note that

$$\begin{aligned}
&(\sum_i c_i X_i)^2 \\
&= (c_1 X_1 + c_2 X_2 + \dots + c_n X_n) \bullet (c_1 X_1 + c_2 X_2 + \dots + c_n X_n) \\
&= \sum_i \sum_j c_i c_j \bullet X_i X_j
\end{aligned}$$

thus

$$\begin{aligned}
&E[(\sum_i c_i X_i)^2] \\
&= E[\sum_i \sum_j c_i c_j \bullet X_i X_j] \\
&= \sum_i \sum_j c_i c_j \bullet E[X_i X_j]
\end{aligned}$$

similarly

$$\begin{aligned}
&E[(\sum_i c_i X_i)]^2 \\
&= [\sum_i c_i \bullet E[X_i]]^2 \\
&= \sum_i \sum_j c_i c_j \bullet E[X_i] E[X_j]
\end{aligned}$$

so

$$\begin{aligned}
& \text{Var}(c^T X_i) \\
&= \sum_i \sum_j c_i c_j \cdot E[X_i X_j] - \sum_i \sum_j c_i c_j \cdot E[X_i] E[X_j] \\
&= \sum_i \sum_j c_i c_j \cdot (E[X_i X_j] - E[X_i] E[X_j]) \\
&= \sum_i \sum_j c_i c_j \cdot \text{Cov}(X_i, X_j) \\
&= c^T \Sigma c
\end{aligned}$$

where Σ is the covariance matrix.

Furthermore, the *first principal component* is the linear combination of x -variables that has maximum variance (among all linear combinations). It accounts for as much variation in the data as possible.

https://stats.stackexchange.com/questions/31177/does-the-variance-of-a-sum-equal-the-sum-of-the-variances?utm_medium=organic&utm_source=google_rich_qa&utm_campaign=google_rich_qa

9. variance $\sigma^2 = (\sum (x - \mu)^2) / N$, captures the spread among the outcomes

10. From the point view of vector, variance σ^2 is the dot product of the data over the number of the data, which means the variance is just the unique length of the whole data distribution. $\sigma^2 = (1/n) XX^T$

11. Similarly, the covariance is defined as $\sigma^2_{XY} = (1/n) XY^T$

If X and Y are correlated, the dot product of XY^T would be non-zero. In contrast, it would be 0, or orthogonal otherwise.

for more details,

<https://imai.princeton.edu/teaching/files/Expectation.pdf>

<https://www.probabilitycourse.com>

Convolution of joint probability distribution:

If X and Y are two jointly continuous random variables and $Z = X + Y$, then

$$f_Z(z) = \int f_{XY}(w, z-w)dw = \int f_{XY}(z-w, w)dw$$

If X and Y are also independent, then

$$\begin{aligned} f_Z(z) &= f_X(z) * f_Y(z) \\ &= \int f_X(w) * f_Y(z-w) * dw \\ &= \int f_Y(w) * f_X(z-w) * dw \end{aligned}$$

for the X and Y be two independent standard normal random variables,
if $X \sim N(\mu_x, \sigma_x^2)$ and $Y \sim N(\mu_y, \sigma_y^2)$ are said to be jointly normal
based on the central limit theorem, if

$$X + Y \sim N(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2 + 2\rho(X, Y)\sigma_x\sigma_y)$$

Moment Generating Function:

$$\begin{aligned} M_X(s) &= E[e^{sX}] \\ &= E[1 + sX + (sX)^2 / 2! + (sX)^3 / 3! + (sX)^4 / 4! + \dots \\ &\quad + (sX)^k / k!] \# \text{ Taylor Series} \\ &= E[\sum (X^k)(s^k) / k!] \\ &= \sum E[X^k] (s^k) / k! \# \text{ expectation linearity} \end{aligned}$$

where $E[X^k]$ is the **k-th moment** of random variable X.

MGF of X gives all moments of X which uniquely determines the distribution.
If two random variables have the same MGF, then they must have the same distribution.

Thus, if you find the MGF of a random variable, you have indeed determined its distribution.

Markov Inequality:

$$P(X \geq a) \leq E[X] / a$$

Chebyshev's Inequality:

$$P(X - \mu \geq a) \leq \text{Var}(X) / a^2$$

Proof: Based on Markov Inequality $P(X \geq a) \leq E[X] / a$,

$$P(X - \mu \geq a) = P((X - \mu)^2 \geq a^2) \leq E[(X - E[X])^2] / a^2 \leq \text{Var}(X) / a^2$$

Chernoff bounds:

Chernoff bounds use of moment generating functions in an essential way to give exponential deviation bounds.

$$P(X - \mu \geq a) \leq E[e^{\lambda(X - \mu)}] e^{-\lambda a}$$

where $E[e^{\lambda(X - \mu)}]$ is the moment generating function of random variable $X - \mu$.

Proof: $P(X - \mu \geq a) = P(e^{\lambda(X - \mu)} \geq e^{-\lambda a}) \leq E[e^{\lambda(X - \mu)}] e^{-\lambda a}$ # Markov Inequality

Hoeffding Lemma:

Let X be a bounded random variable with $X \in [a, b]$, then

$$E[e^{\lambda(X - \mu)}] \leq \exp(\lambda^2(b - a)^2 / 8) \text{ for all } \lambda \in \mathbb{R}$$

Hoeffding inequality:

Let X_1, \dots, X_n be independent bounded random variables with $X_i \in [a, b]$ for all i , where $-\infty < a \leq b < \infty$, then

$$P((1/n) \sum (X_i - \mu_i) \geq t) \leq \exp(-2nt^2 / (b - a)^2)$$

for all $t \geq 0$.

Proof:

$$\begin{aligned} P((1/n) \sum (X_i - \mu_i) \geq t) &= P(\sum (X_i - \mu_i) \geq nt) \\ &\leq E[e^{\lambda \sum (X_i - \mu_i)}] e^{-\lambda nt} \quad \# \text{ Chernoff} \end{aligned}$$

bounds

$$\leq \exp(-2nt^2 / (b - a)^2) \quad \# \text{ Hoeffding}$$

Lemma

<http://cs229.stanford.edu/extr-notes/hoeffding.pdf>

Temperature definition in thermodynamics theory:

Temperature is the change of the expected energy in a thermodynamic equilibrium system, which is the sum of probability of each microstate times the energy required by the state,

that increasing 1 bit of entropy to the system, which is the negative sum of the probability of each microstate times the log of it.

$$T = dE / dH$$

1st law of thermodynamics: energy remains the same among systems

2nd law of thermodynamics: the entropy of the whole system increases.

we introduced 2 energies here: E1 and E2:
first thing first, from the 1st law of thermodynamics:
 $dE1 + dE2 = 0$

by the definition of temperature,
 $dE1 = T1 * dH1$, $dE2 = T2 * dH2$

and since the entropy of the 2 systems each of which owns its own energy increases, we have

$$dH1 + dH2 > 0$$

thus,

$$dH1 - (T1 / T2) * dH1 > 0$$

which means,

**If system didn't reach to the thermodynamic equilibrium,
which means the temperature is not equal, entropy increases.**

the thermodynamic equilibrium is rarely happened.

Euler Number e:

1. $e = \lim (1 + 1/n)^n = \int (1 / n!) dn$,
it could be explained as, in an other way, the Taylor Series of function e^x at the point $x = 1$

<https://www.youtube.com/watch?v=3d6DsJIBzJ4>

$$2. e^x = \lim (1 + x/n)^n = \int (x / n!) dn$$

3. e is always being used in the equations as to represent the factorial facets like permutations etc.

Definition of ln:

$$\ln(x) = \int 1/x \, dx$$

Computing Log-Sum-Exp

Imagine that we have a set of N values and we want to compute the quantity

$$z = \log \sum_i \exp \{ x_i \}$$

For extremely large or extremely negative x_i , z would be resulted in underflow or overflow. With some transformations we have,

$$\exp \{ z \} = \sum_i \exp \{ x_i \}$$

$$\exp \{ -\alpha \} \cdot \exp \{ z \} = \exp \{ -\alpha \} \cdot \sum_i \exp \{ x_i \}$$

$$\exp \{ z - \alpha \} = \sum_i \exp \{ x_i - \alpha \}$$

$$z - \alpha = \log \sum_i \exp \{ x_i - \alpha \}$$

$$z = \alpha + \log \sum_i \exp \{ x_i - \alpha \}$$

A typical choice of α is

$$\alpha = \operatorname{argmax}_i x_i$$

which ensures that the largest value you exponentiate will be zero.

z can be roughly approximated by α , that is $\operatorname{argmax}_i x_i$. This approximation is based on the idea that $\exp \{ x_k \}$ is insignificant for any x_k that is noticeably less than $\operatorname{argmax}_i x_i$.

Exponentially Weighted Average

Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1 - \beta) \theta_t$$

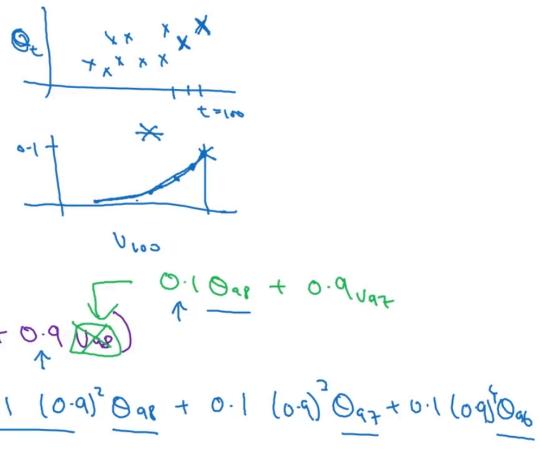
$$v_{100} = 0.9 v_{99} + 0.1 \theta_{100}$$

$$v_{99} = 0.9 v_{98} + 0.1 \theta_{99}$$

$$v_{98} = 0.9 v_{97} + 0.1 \theta_{98}$$

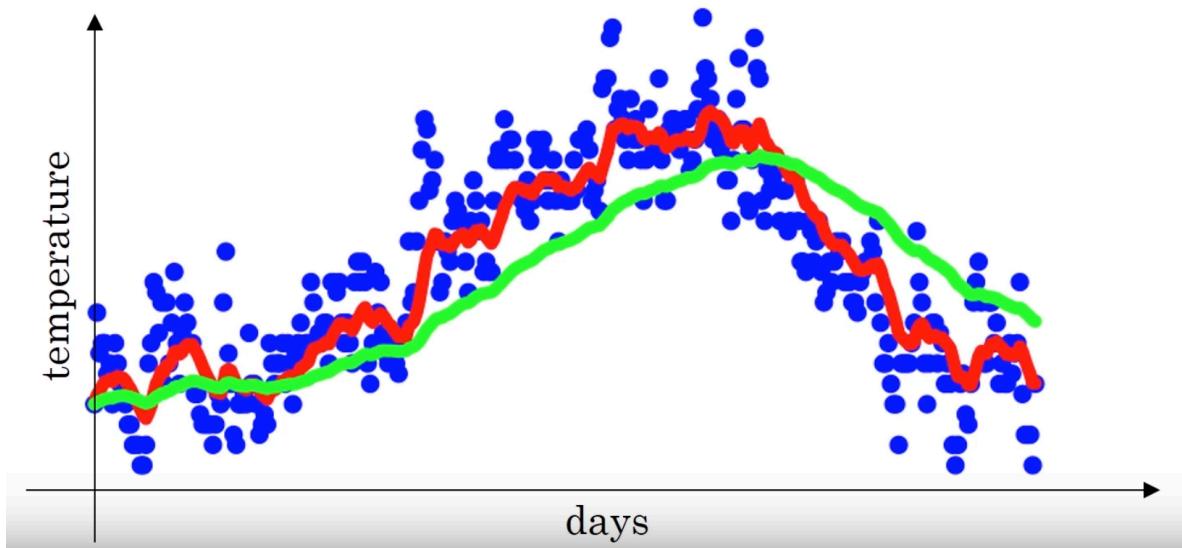
...

$$\begin{aligned} v_{100} &= 0.1 \theta_{100} + 0.9 v_{99} \\ &= 0.1 \theta_{100} + 0.9 (0.1 \theta_{99} + 0.9 v_{98}) \\ &= 0.1 \theta_{100} + 0.9 \cdot 0.1 \theta_{99} + 0.9^2 v_{98} \\ &\quad + \dots \end{aligned}$$



which is exponentially average $1 / (1 - \beta)$ historical data to compute the current data.

The larger the β , the smoother the weighed average. As shown below, $\beta = 0.9$ generates the red line, while $\beta = 0.98$ generates the green line.



<https://www.youtube.com/watch?v=NxTFIzBjS-4>

Gaussian distribution:

$$1. N(\mu, \sigma^2) = \mu + N(0, \sigma^2)$$

2. for the linear transformation $y = w^T x$ with some gaussian noise $\epsilon \sim N(0, \sigma^2)$

$$y = w^T x + \epsilon$$

recall that for a normally distributed random variable with mean m and variance h : $t \sim N(m, h)$, then a linear transformation of t given by $at + b$ is also a normally distributed random variable with mean $a * m + b$, and variance $h * a^2$.

thus, y is the linear transformation of random variable ϵ with mean $w^T x$, and variance σ^2 , which is

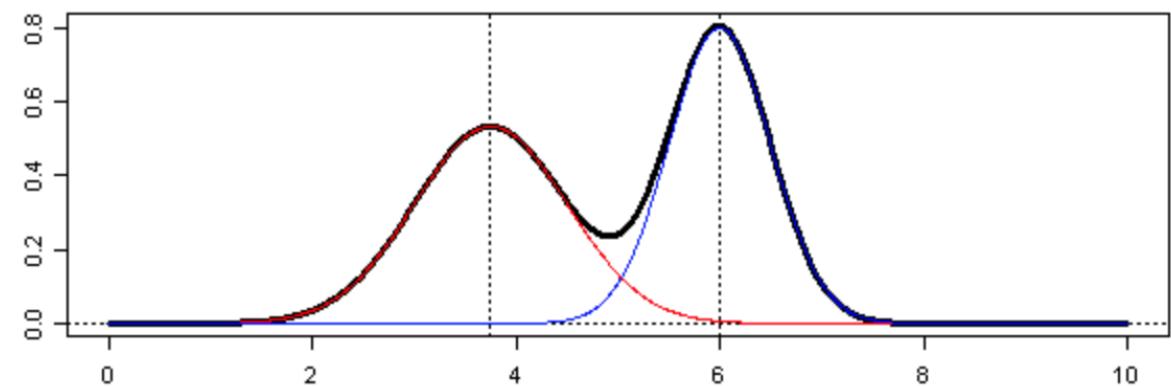
$$y | x \sim N(w^T x, \sigma^2)$$

Furthermore, if given that $x \sim N(0, I)$, thus

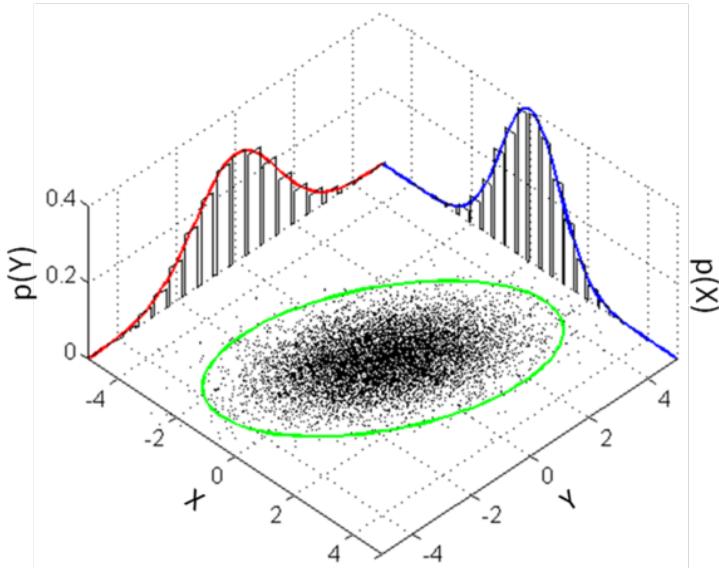
$$\begin{aligned} & \text{Var}(y | x) \\ &= \text{Var}(w^T x + \epsilon) \\ &= w^T \cdot \text{Var}(x) \cdot w + \text{Var}(\epsilon) \\ &= w^T w + \sigma^2 \end{aligned}$$

3. Gaussian Mixture Model vs. Multivariate Gaussian Model

$$\text{GMM} = \sum_i N(x; \Phi_i)$$



$$\text{MGM} = \mathbf{N}(X_{1:N}; \Phi)$$



Bernoulli distribution function:

$$1. P(x; \theta) = \theta^x * (1 - \theta)^{1-x}$$

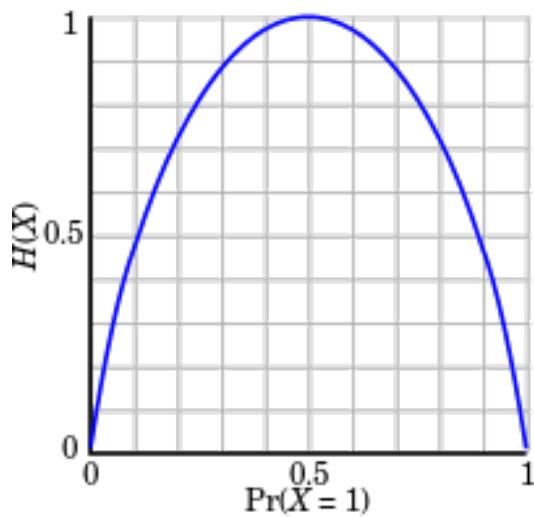
for x belongs to $\{0, 1\}$, the common estimator (**sample mean**) for θ is $(1/m) \sum_i x_i$

$$2. E(x) = \sum_i x_i * p_i = (1 * \theta^1 * (1 - \theta)^{1-1}) + (0 * \theta^0 * (1 - \theta)^{1-0}) = \theta$$

for x belongs to $\{0, 1\}$

3. the entropy for a Bernoulli variable X

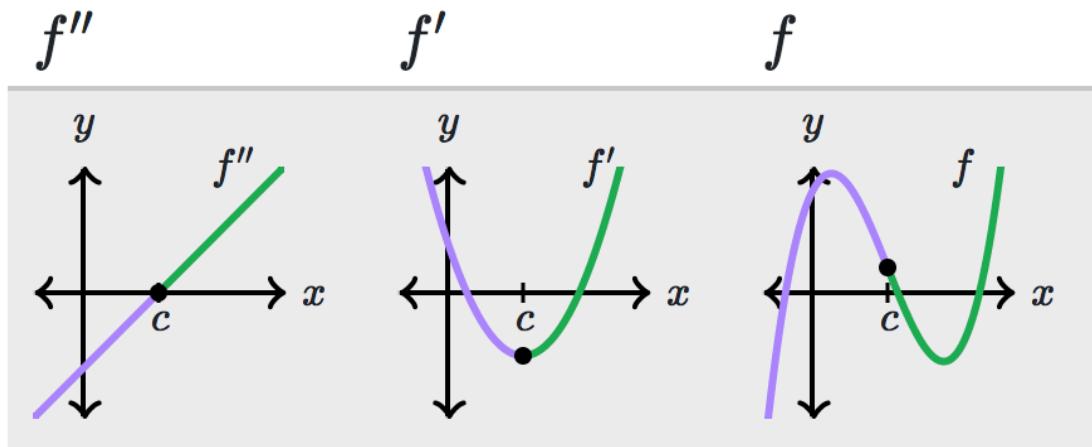
$$H(x) = -\sum P(x; \theta) \log(P(x; \theta)) = -((1 - \theta) \log(1 - \theta) + \theta \log \theta)$$



How f'' informs about the concavity of f

When the second derivative f'' is positive, that means the first derivative f' is increasing, which means that f is concave up.

Similarly, a negative f'' means f' is decreasing and f is concave down.



Taylor Series:

1. <https://www.youtube.com/watch?v=3d6DsJIBzJ4>
2. the basic intuition is it transform the derivative information at a point to

the output information near that point. Precisely, if we want to approximate the curve of $\cos(\theta)$ at the point 0 with the polynomial function controlled by some parameters \mathbf{C}_N ,

let's say

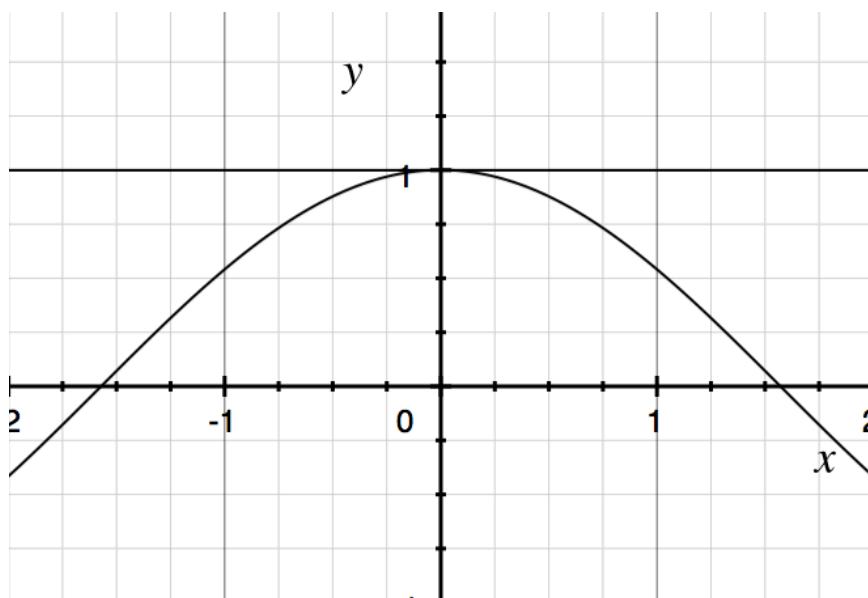
$$f(\theta) = C_0 + C_1 \theta + C_2 \theta^2 + C_3 \theta^3 + \dots + C_N \theta^N$$

The choice of \mathbf{C}_N is determined by the derivative of $\cos(\theta)$ and $f(\theta)$ at the point 0 .

Let $N = 0$, $f(\theta) = C_0$

We approximate $\cos(\theta)$ at the point 0 as a constant line. Let $\cos(0) = f(0)$, then $C_0 = 1$.

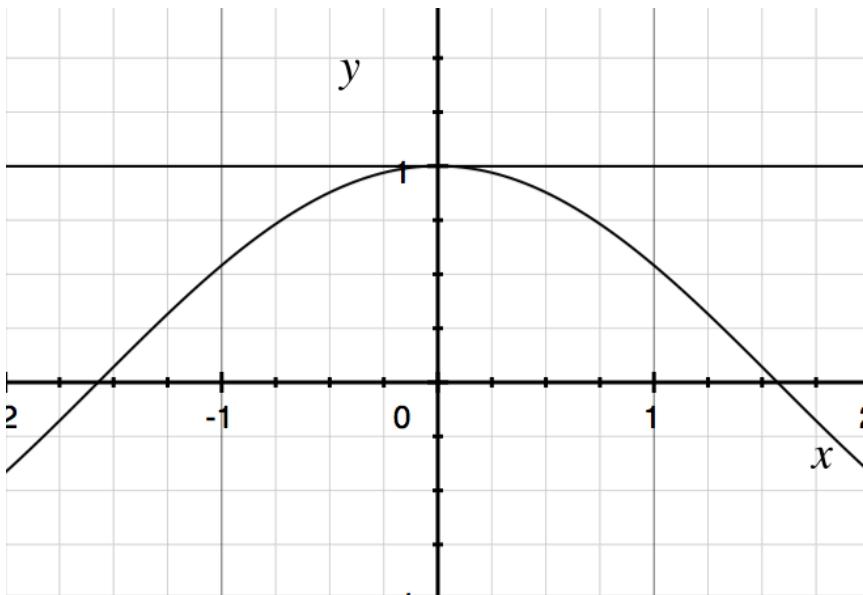
$$f(\theta) = 1$$



Let $N = 1$, $f(\theta) = 1 + C_1 \theta$

We approximate $\cos(\theta)$ at the point 0 as a line with slope C_1 . Let $\partial \cos(\theta) / \partial \theta |_{\theta=0} = \partial f(\theta) / \partial \theta |_{\theta=0}$, then $C_1 = 0$.

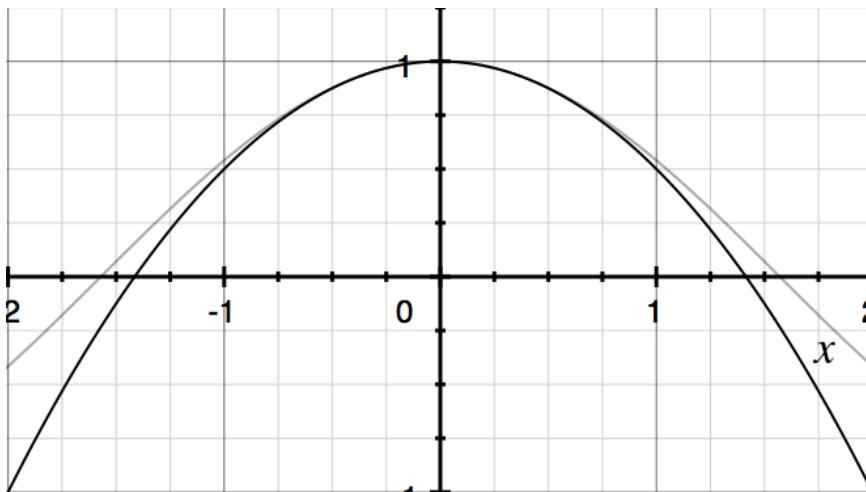
$$f(\theta) = 1$$



Let $N = 2$, $f(\theta) = 1 + C_2 \theta^2$

We approximate $\cos(\theta)$ at the point 0 as a quadratic curve. Let $\partial^2 \cos(\theta) / \partial \theta^2$
 $|_{\theta=0} = \partial^2 f(\theta) / \partial \theta^2 |_{\theta=0}$, then $C_2 = -1/2$.

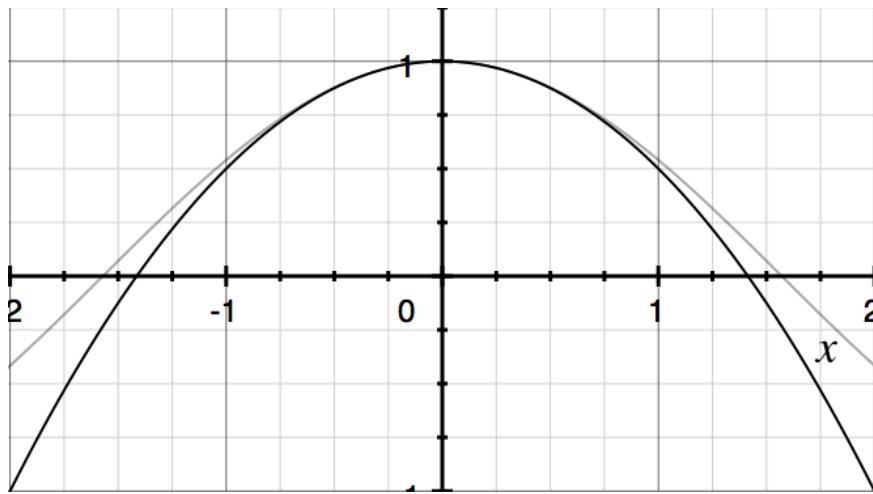
$$f(\theta) = 1 - (1/2)\theta^2$$



Let $N = 3$, $f(\theta) = 1 - (1/2)\theta^2 + C_3 \theta^3$

Let $\partial^3 \cos(\theta) / \partial \theta^3 |_{\theta=0} = \partial^3 f(\theta) / \partial \theta^3 |_{\theta=0}$, then $C_3 = 0$.

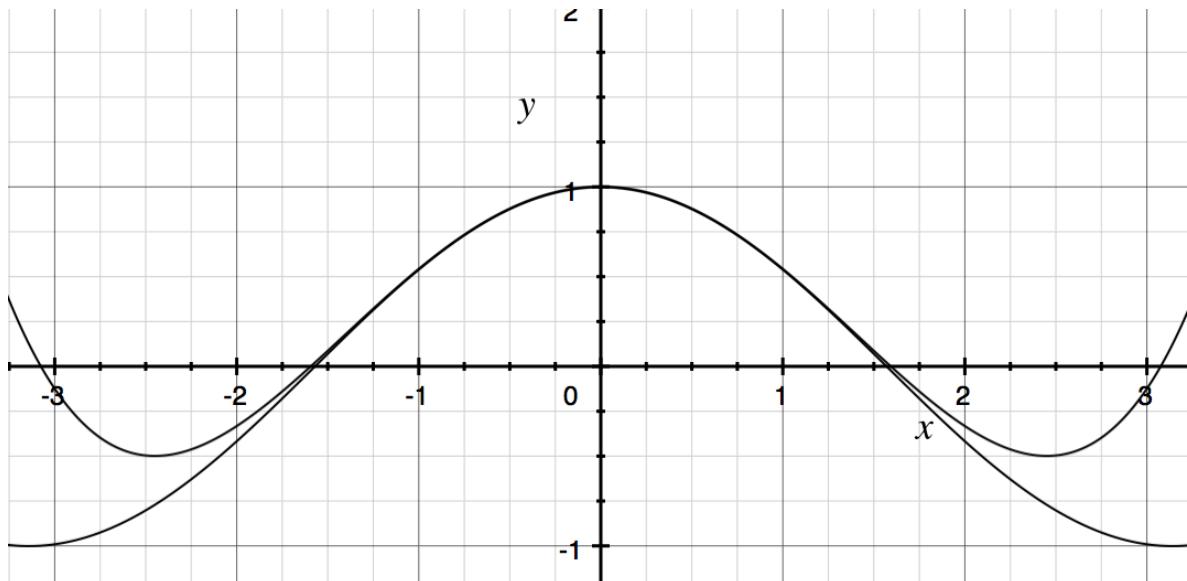
$$f(\theta) = 1 - (1/2)\theta^2$$



Let $N = 4$, $f(\theta) = 1 - 0.5\theta^2 + C_4 \theta^4$

Let $\frac{\partial^4 \cos(\theta)}{\partial \theta^4} \Big|_{\theta=0} = \frac{\partial^4 f(\theta)}{\partial \theta^4} \Big|_{\theta=0}$, then $C_4 = 4!$.

$$f(\theta) = 1 - (1/2)\theta^2 + (1/4!) \theta^4$$



...

The larger the exponent, the deeper the derivative matching, the more precise approximation near the point.

More generally, if $f(x)$ is the approximation of the function $g(x)$ at the point θ , with Taylor Series we have,

$$f(x) = g(\emptyset) + (\partial g(x) / \partial x|_{x=\emptyset} (x - \emptyset)^1 / 1!) + (\partial^2 g(x) / \partial x^2|_{x=\emptyset} (x - \emptyset)^2 / 2!) + \dots + (\partial^n g(x) / \partial x^n|_{x=\emptyset} (x - \emptyset)^n / n!)$$

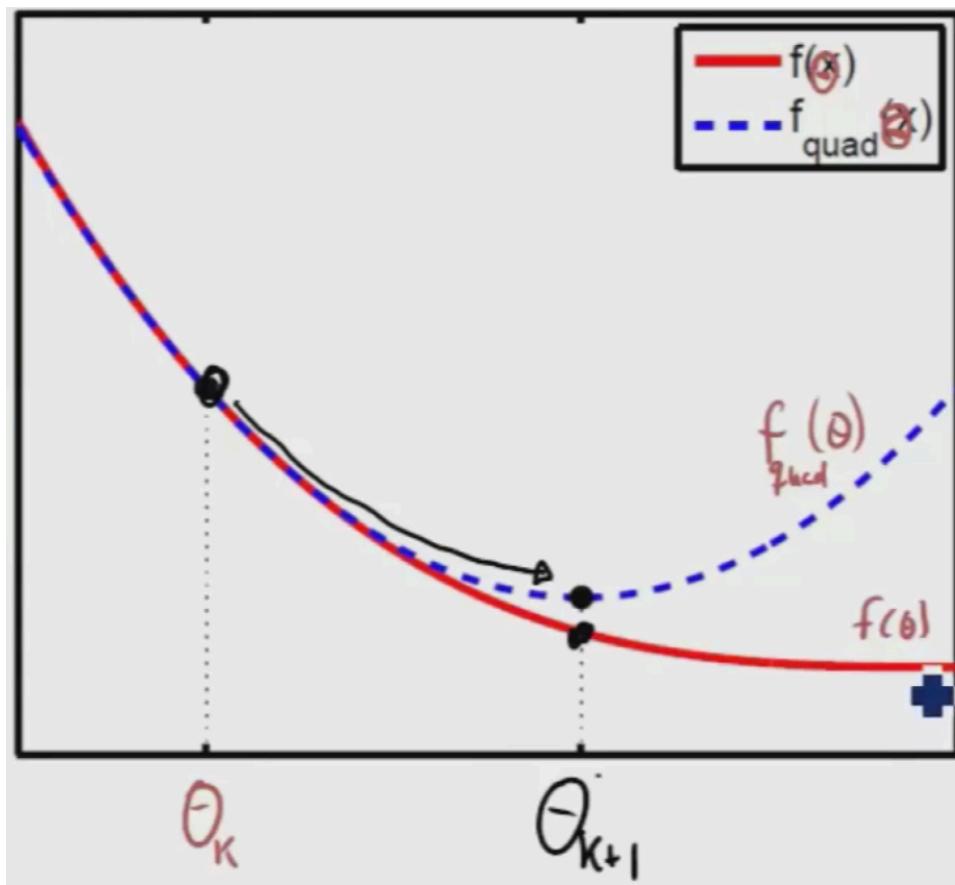
check out <https://www.youtube.com/watch?v=ClFrIgOPpnM> for multivariate Taylor Series.

Newton's Method for Optimization

Making the **2nd-order Taylor Series** approximation of the cost function $f(\theta)$ around θ_k , we will get a quadratic approximation of $f(\theta)$ near θ_k ,

$$f_{\text{quad}}(\theta) = f(\theta_k) + (\theta - \theta_k)^T g_k + (1/2)(\theta - \theta_k)^T H_k (\theta - \theta_k)$$

where g_k and H_k are the gradient and hessian of $f(\theta)$ at the point θ .



So the next point θ_{k+1} for the steepest step of the cost function f can be calculated by finding the **zero gradient** point of the quadratic approximation function, $\nabla f_{\text{quad}} = \mathbf{0}$, thus,

$$\theta_{k+1} = \theta_k - H_k^{-1} \cdot g_k$$

And that is the brilliant **Newton's Method**.

Definition of Residual:

In **regression analysis**, the difference between the observed value of the dependent variable (y) and the predicted value (\hat{y}) is called the **residual** (e). Each data point has one residual.

Definition of Gamma function:

1. $\Gamma(n + 1) = n! = \int e^{-t} * t^n * dt [1] = (z - 1)!$
2. using the Taylor Series, we can approximate $e^{-t} * t^n \approx (n^n / e^n) * e^{(-\varepsilon^2 / 2n)} [2]$, and taking [2] into the integration [1], we get the Stirling Approximation Equation.
3. the Gamma prior is conjugated to the Poisson or Gaussian distribution.

Definition of Beta function:

1. Beta function is used for modeling the prior of Bernoulli distribution in **Bayesian learning procedure**.

$$\begin{aligned}
f(x; \alpha, \beta) &= \text{constant} \cdot x^{\alpha-1} (1-x)^{\beta-1} \\
&= \frac{x^{\alpha-1} (1-x)^{\beta-1}}{\int_0^1 u^{\alpha-1} (1-u)^{\beta-1} du} \\
&= \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1} \\
&= \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1}
\end{aligned}$$

2. as α, β being determined, the posterior probability can be calculated with the conjugate of prior probability and likelihood, instead of the hard work of the marginal probability calculation using integral which might be a NP problem in high dimensional space.
3. Beta prior is conjugated to the Bernoulli distribution.

Representer Theorem

Suppose the regularized risk is defined as

$$J_\lambda(\theta) = (1/m) \sum_i \text{Loss}(\theta^T x_i, y_i) + (\lambda/2) \|\theta\|^2$$

Then there is a minimizer of the risk (0 gradient) that can be written as

$$\theta = \sum_i \alpha_i x_i$$

where $\alpha_i = -\text{Loss}'(\theta^T x_i, y_i) / \lambda$.

Linear Transformation

1. The **prime** form of a linear transformation

$$\mathbf{y} = \boldsymbol{\theta}^T \mathbf{x}$$

$\boldsymbol{\theta}$ is perpendicular to \mathbf{y} .

2. Based on the representer theorem, the linear transformation can be always written as the **dual** form

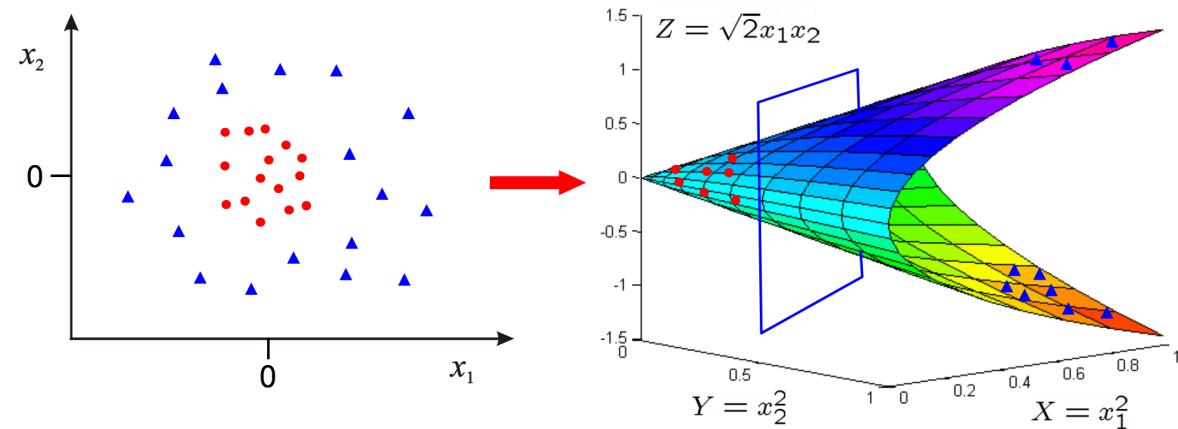
$$\mathbf{y} = \sum_i \alpha_i \mathbf{x}^T \mathbf{x}_i$$

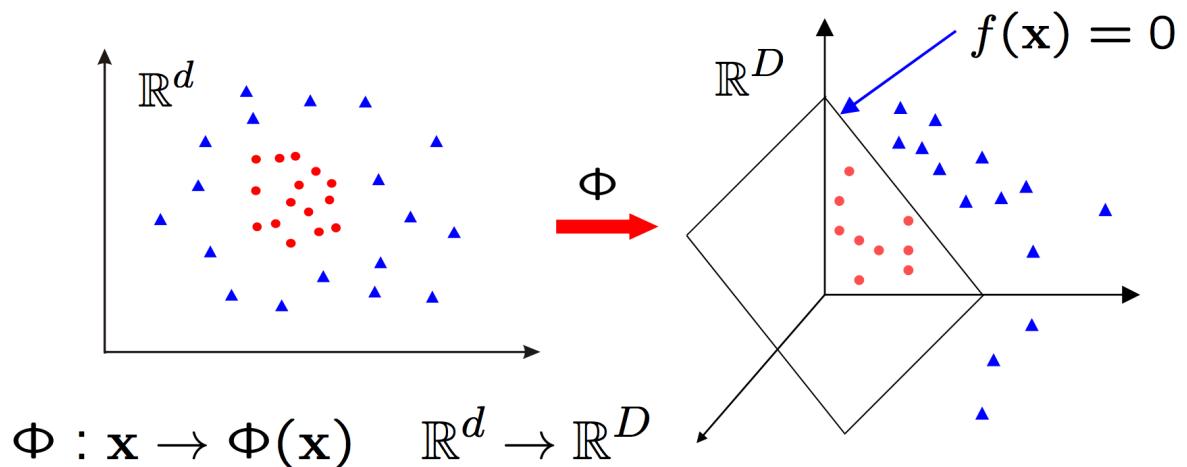
That is, in any learning algorithm, minimize directly over $\boldsymbol{\alpha} \in \mathbb{R}^m$

Kernel Trick

1. Sometimes data are not linearly separable in low dimensional space, it can be solved by mapping the data to higher dimension.

$$\Phi : \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \rightarrow \begin{pmatrix} x_1^2 \\ x_2^2 \\ \sqrt{2}x_1x_2 \end{pmatrix} \quad \mathbb{R}^2 \rightarrow \mathbb{R}^3$$





the **Inner Product** of the mapping function Φ is called **Kernel**.

$$K(\mathbf{x}, \mathbf{x}) = \Phi(\mathbf{x})^T \Phi(\mathbf{x}).$$

Then we can rewrite the regularized risk as

$$\begin{aligned} J_\lambda(\alpha) &= (1/m) \sum_i \text{Loss}(\sum_j \alpha_j K(\mathbf{x}_i, \mathbf{x}_j)) + (\lambda/2) \sum_{ij} \alpha_i \alpha_j \\ K(\mathbf{x}_i, \mathbf{x}_j) \end{aligned}$$

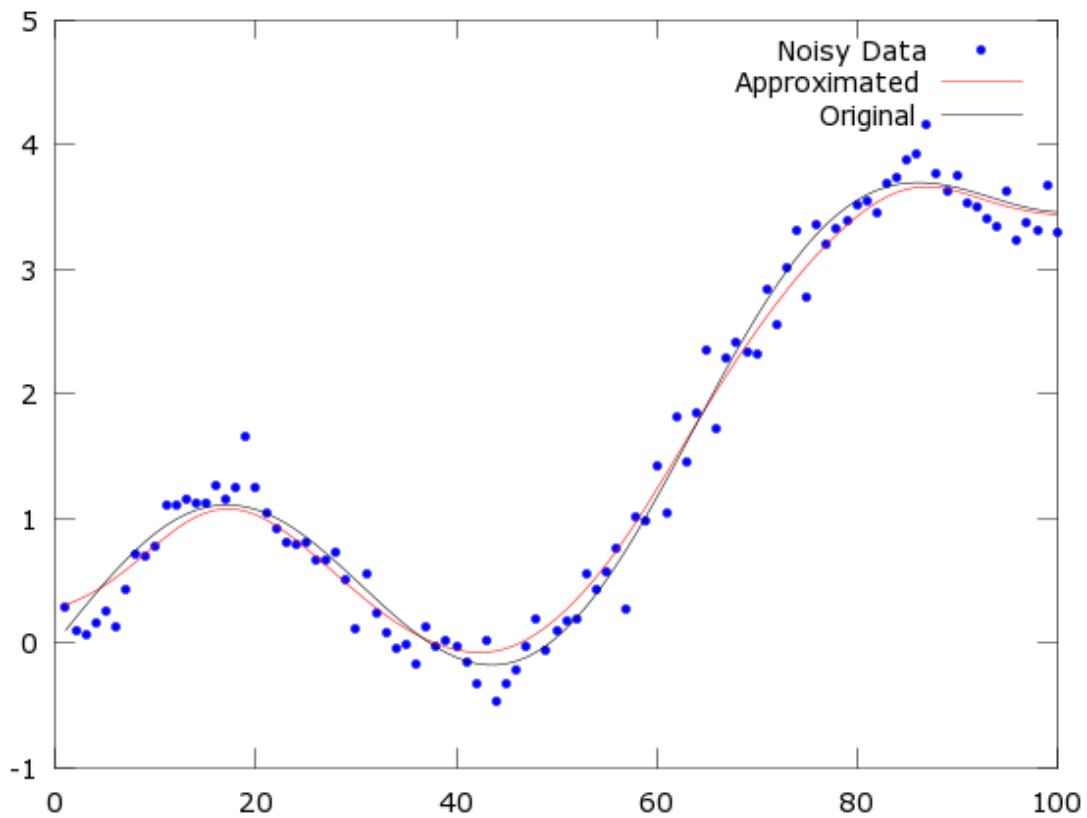
2. Kernel Trick

The kernel function $K(\mathbf{x}, \mathbf{z})$ may be very inexpensive to calculate, even though $\Phi(\mathbf{x})$ itself may be very expensive to calculate (perhaps because it is an extremely high dimensional vector). In such settings, we can learn higher dimensional features spaces, without actually having to map the points into the high dimensional space, say calculate $\Phi(\mathbf{x})$.

For example, the Gaussian or Radial Basis Function (RBF) kernel is applicable to data in any dimension.

$$K(\mathbf{x}, \mathbf{x}_0) = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2 / 2\sigma^2) \text{ for } \sigma > 0$$

Gaussian Kernel Regression:



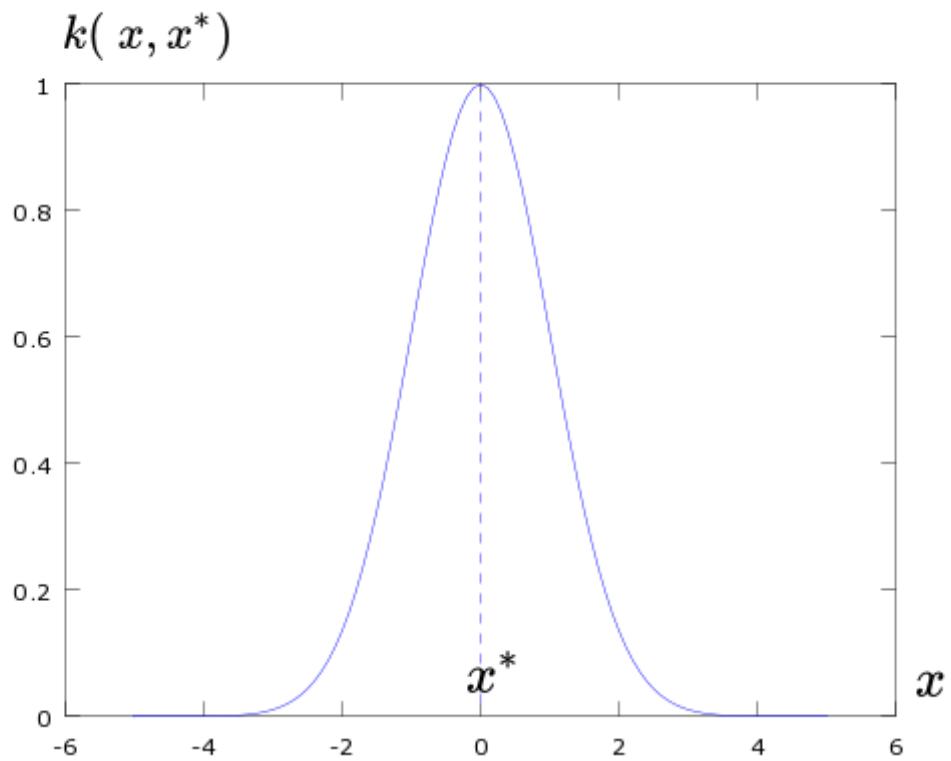
Gaussian Kernel Regression is a technique for **non-linear** regression. Instead of the **Iterative Learning** method (such as gradient decent in DNN), Gaussian Kernel Regression approximate a function using **Weighted Average** method, which is the prediction **y-hat** could be computed by **the history of y with some chosen weights** introduced.

$$\bar{y} = \frac{\sum_{i=1}^m (w_i y_i)}{\sum_{i=1}^m w_i}$$

So how to choose the weight ?

The intuition here is that if the history point is closer to the query point, it should be considered to be important to the prediction (a higher weight). Likewise, if the history point is far away from the query point, it should be considered to be negligible to the prediction (a lower weight).

The gaussian function has exactly the behavior we want for computing the weight values.



$$k(x, x^*) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-x^*)^2}{2\sigma^2}}$$

The denominator will cancel out when compute the weighted average, thus

$$k(x, x^*) = e^{-\frac{(x-x^*)^2}{2\sigma^2}}$$

which is called the **Gaussian Kernel**.

<http://mccormickml.com/2014/02/26/kernel-regression/>

Gaussian Process:

A GP is a Gaussian distribution over functions $f(x_i)$.

$$f(x) \sim N(e(x), k(x, x')) + \epsilon,$$

$$e(x) = E[f(x)],$$

$$\epsilon \sim N(0, \sigma^2),$$

$$k(x, x') = E[(f(x) - e(x))(f(x') - e(x'))^T] + \sigma^2 I,$$

$k(x, x')$ is called the **kernel function** which finds out the covariance of $f(x)$ and $f(x')$ as given x and x' , plus the variance of the noise.

There are different ways to choose kernel function for different problems, here choose

$$k(x, x') = \exp(-(1/2\ell^2)(x - x')^2)$$

for GP.

Prediction:

As new data comes in, the joint probability of f, f^* is

$$\begin{pmatrix} f \\ f_* \end{pmatrix} \sim \mathcal{N}\left(\begin{pmatrix} \mu \\ \mu_* \end{pmatrix}, \begin{pmatrix} K & K_* \\ K_*^T & K_{**} \end{pmatrix}\right)$$

Then we can get $p(f_* | f)$ with the formula of multivariate gaussian conditional probability.

$$\begin{aligned} p(f_* | X_*, X, f) &= \mathcal{N}(f_* | \mu_*, \Sigma_*) \\ \mu_* &= \mu(X_*) + K_*^T K^{-1} (f - \mu(X)) \\ \Sigma_* &= K_{**} - K_*^T K^{-1} K_* \end{aligned}$$

Gaussian Process is equivalent to **Ridge Regression**, which means introducing the regularization is introducing noise into the hypotheses actually.

Moreover, if you take a one-layer hidden neural network, and just make the number of neurons goes infinity, this NN actually is a Gaussian Process Model.

<https://www.youtube.com/watch?v=MfHKW5z-OOA&index=9&list=PLE6Wd9FR--EdyJ5lbFI8UuGjecvVw66F6>
<https://www.youtube.com/watch?v=ewJ3AxKcI0g>

Exponential Family:

The Exponential Family of distributions of y as those having the form:

$$p(y; \eta) = b(y) \exp\{ \eta^T T(y) - a(\eta) \}$$

1. η is called the natural parameter (also called the canonical parameter) of the distribution
2. $T(y)$ is the sufficient statistic ($T(y)$ is accurate enough to describe the distribution. $T(y) = y$ in general)
3. $a(\eta)$ is the log *partition function*. The quantity $e^{-a(\eta)}$ essentially plays the role of normalization constant, that makes sure the distribution $p(y; \eta)$ sums / integrates over y to 1.

In thermodynamics, partition function is calculated by minimizing the entropy over P_i constrained by average energy, average momentum etc., with the LaGrange multiplier. It is the most important part for the Canonical ensemble.

$$4. -\nabla a(\eta) = E[T(y)]$$

Generalized Linear Model (GLM) :

Given a data set X (imagine the pixels input from training images), we want to understand the expectation of the (conditional) distribution of $Y | X$, in particular we would like to use the linear predictor $\eta = \beta^T x$. In the linear model, we use η as the predictor for $E[Y | X]$

A generalization of this is the Generalized Linear Model (GLM):

$$\eta = \beta^T x = g(E[Y | X])$$

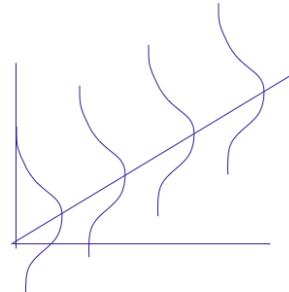
where $g(\bullet)$ is some link function (it links the **expectation** of the distribution of Y to the linear expression).

Generalised linear models – components

- **Distribution** of Y_i 's:

linear models $Y_i \sim N(\mu_i, \sigma^2)$

GLMs $Y_i \sim \text{exponential family}$



- **Linear predictor** = function of the covariates:

linear models $\eta_i = \alpha + \beta x_i$

GLMs eg $\eta_i = \alpha + \beta x_i + \gamma x_i^2$

- **Link function** = connection between the linear predictor and $\mu_i = E(Y_i)$:

linear models $\eta_i = \mu_i \Rightarrow \mu_i = \alpha + \beta x_i$

GLMs eg $\eta_i = \ln \mu_i \Rightarrow \mu_i = e^{\alpha + \beta x_i + \gamma x_i^2}$

Gradient:

1. The sum of all the partial derivative vectors of a function. It is a vector that always perpendicular to the contour line of the function.

$$\operatorname{div} \mathbf{F} = \nabla \cdot \mathbf{F} = \left(\frac{\partial}{\partial x}, \frac{\partial}{\partial y}, \frac{\partial}{\partial z} \right) \cdot (U, V, W) = \frac{\partial U}{\partial x} + \frac{\partial V}{\partial y} + \frac{\partial W}{\partial z}.$$

From the definition of derivative, the $\partial F / \partial x$ represent how much F change respective to the slitty change of x. It can be also interpreted as, the value of vector ∂F divided by vector ∂x , the result is also a vector.

2. The derivative of the gradient, also the second order derivative $\partial^2 F / \partial x^2$, represent the **sharpness** of the curve.

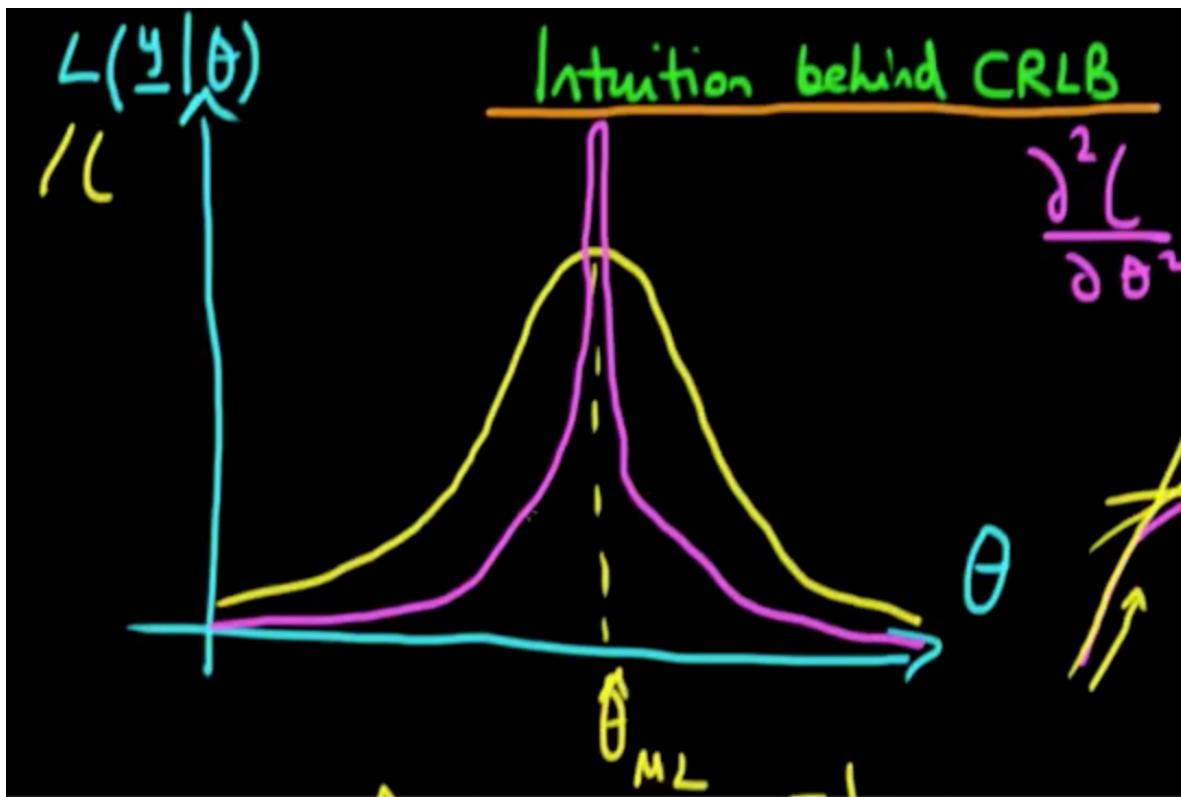
The higher the $\partial^2 F / \partial x^2$, the sharper the curve, thus the low variance of the distribution in probability distribution function (PDF).

The Relationship between the Curvature (2nd derivative, or Information matrix I) of Log-likelihood function and the Variance: Cramer-Rao Lower Bound

$$\operatorname{Var}(\theta) \geq I^{-1}(\theta) = [\partial^2 L(\theta) / \partial \theta^2]^{-1}$$

Think of it in geometry, a gaussian distribution for instance,

more narrow the bell -> lower variance the distribution has -> larger the curvature -> less uncertainty



Furthermore, the Fisher Information Matrix is the negative Hessian Matrix:

$$I(\theta) = -H(\theta) = -\frac{\partial^2 L(\theta)}{\partial \theta_i \partial \theta_j}, \quad 1 \leq i, j \leq p$$

Maximum Likelihood Estimation (MLE)

Maximum Likelihood Estimation is a **frequentist (discriminative) statistics** approach that can be used to derive parameter (point) estimators in a parametric estimation problem.

$$\begin{aligned}
 \theta_{ML} &= \operatorname{argmax}_{\theta} L(X; \theta) \\
 &= \operatorname{argmax}_{\theta} \prod_i L(x_i; \theta) \\
 &= \operatorname{argmax}_{\theta} \sum_i \log L(x_i; \theta) \\
 &= \operatorname{argmax}_{\theta} (1/N) \sum_i \log L(x_i; \theta) \quad \# \text{ Monte Carlo estimate} \\
 &\text{of the true cross entropy} \\
 &= \operatorname{argmax}_{\theta} \sum_i p(x_i) \log L(x_i; \theta) \\
 &= \operatorname{argmin}_{\theta} -\sum_i p(x_i) \log L(x_i; \theta) \quad \# \text{ cross entropy}
 \end{aligned}$$

We can thus see maximum likelihood as an attempt to make the model distribution match the empirical distribution . Any loss consisting of a negative log-likelihood is a cross entropy between the empirical distribution defined by the training set and the model. We often phrase it as minimizing a cost function.

More generally, MLE can be used to estimate a conditional probability $P(Y | X ; \theta)$ in order to predict y given x for **classification** and **regression** problems.

$$\theta_{ML} = \operatorname{argmax}_{\theta} \sum_i \log L(y_i | x_i; \theta)$$

For example, mean squared error is the cross-entropy between the empirical distribution and a gaussian model. That is each data point is seen as gaussian distributed.

Proof:

for the linear transformation $y = w^T x$ with some gaussian noise $\epsilon \sim N(0, \sigma^2)$

$$y = w^T x + \epsilon$$

recall that for a normally distributed random variable with mean m and variance h : $t \sim N(m, h)$, then a linear transformation of t given by $at + b$ is also a normally distributed random variable with mean $a * m + b$, and variance $h * a^2$.

y is the linear transformation of random variable ϵ with mean $w^T x$, and variance σ^2 , which is

$$y | x \sim N(w^T x, \sigma^2)$$

thus

$$P(y | x) \sim N(y; w^T x, \sigma^2)$$

the estimator of $f(y | x)$ using MLE is

$$\begin{aligned} w_{ML} &= \operatorname{argmax}_w (1/N) \sum_i \log p(y_i | x_i; w) \quad \# \text{cross entropy} \\ &\quad \text{between the empirical distribution and a gaussian model} \\ &= \operatorname{argmax}_w -\log\sigma - (1/2) \log(2\pi) - (1/N) \sum_i [\|w^T x_i - y_i\|^2 / 2\sigma^2] \end{aligned}$$

here we assume that the variance is fixed to some constant σ^2 chosen by the user.

Maximizing the log-likelihood is equivalent to minimizing $(1/N) \sum_i \|w^T x_i - y_i\|^2$, which is equivalent to the **MSE**

$$MSE_{train} = (1/N) \sum_i \|w^T x_i - y_i\|^2$$

The MLE is asymptotically normal

That is, as **training size** goes infinity, we have

$$\theta_{ML} \sim N(\theta, I^{-1})$$

where I is the negative of the Hessian Matrix (Fisher Information Matrix) of the log-likelihood function. According to the Cramer-Rao Lower Bound (see above), the MLE attains the lowest variance, which is I^{-1} , and is often considered the preferred estimator to use for machine learning.

Maximum A Posteriori (MAP):

Maximum A Posteriori is a **bayesian (generative) statistics** approach that consider all possible values of θ when making a prediction.

$$\theta_{MAP} = \operatorname{argmax}_{\theta} \log p(y | x, \theta) + \log p(\theta)$$

MAP has the same formulas as for the MLE (maximum likelihood estimator) for θ , except that MLE look at the true parameter θ as a unknown parameter but not a random variable, whereas MAP represents it as a random variable. So it put the prior of θ into account. In practical applications, a common choice for prior $p(\theta)$ is to assume that $\theta \sim N(0, \tau I^2)$ with the highest entropy to reflect a high degree of uncertainty in the value of θ before any observation, or **Beta** function, or using the **Monte Carlo** to sample the prior distribution.

The functionality of **prior to the MAP** is sort of like the **regularization to the MLE**, that both can be used to **prevent overfitting**.

For example, if consider a linear regression model with a gaussian prior on the weights w . If this prior is given by

$$p(w) = N(w; 0, \tau I^2),$$

then

$$\log p(w) \propto (1/\tau) w^T w$$

which is a **L2-Norm** regularizer.

Logistic Regression with MLE

1. Deriving the Log-likelihood

Since the pdf of y is a Bernoulli distribution $y | x \sim \text{Bernoulli}(\theta)$, we assume the conditional likelihood of all data are complied with:

$$p(y | x) = \text{Bernoulli}(\theta) = \theta^y \cdot (1 - \theta)^{1-y}$$

where θ is the expectation of the Bernoulli distribution.

In the form of exponential family $y | x \sim b(y) \exp\{\eta^T T(y) - a(\eta)\}$, the distribution of y is:

$$p(y | x) = \exp\{\log(\theta / (1 - \theta)) \cdot y + \log(1 - \theta)\}$$

where

1. $\eta = \log(\theta / (1 - \theta))$
2. $T(y) = y$
3. $a(\eta) = -\log(1 - \theta) = \log(1 + e^\eta)$
4. $b(y) = 1$

then,

$$\theta = 1 / (1 + e^{-\eta})$$

thus the **link function** is a **logistic function**.

By definition of Generalized Linear Model, η is the **linear predictor** to the expectation θ which is $w^T x$, then we have

$$\theta = 1 / (1 + \exp \{-w^T x - b\})$$

Taking back to the $p(y | x)$ we have

$$p(y | x; w) = [1 / (1 + \exp \{-w^T x - b\})]^y \cdot [1 - 1 / (1 + \exp \{-w^T x - b\})]^{1-y}$$

2. Maximizing Likelihood Estimator:

$$\begin{aligned} w_{ML} &= \operatorname{argmax}_w \log L(y | x; w) \\ &= \operatorname{argmax}_w \log \prod_i [1 / (1 + \exp \{-w^T x_i - b\})]^y_i \cdot [1 - 1 / (1 + \exp \{-w^T x_i - b\})]^{1-y_i} \\ &= \operatorname{argmax}_w \sum [y_i * \log(1 / (1 + \exp \{-w^T x_i - b\}))] + [(1 - y_i) * \log(1 - 1 / (1 + \exp \{-w^T x_i - b\}))] \quad \# \text{cross-entropy, convex} \end{aligned}$$

the derivative of log-likelihood function respect to w is

$$dw = \sum_i x_i^T (1 / (1 + \exp \{-w^T x_i - b\}))$$

then update w by dw :

$$w = w + \xi \cdot dw$$

LaGrange Multiplier:

the basic idea is that as maximizing the original function with the constraint function, both of the contour line of these 2 functions at the maximized value should have the same gradient. thus, $\nabla f = \lambda \nabla g$. In other words, λ is the derivative of function f

respect to the function g.

1. adding the constraints to the target with LaGrange parameter
2. take the partial derivatives of each variable, $\frac{\partial f}{\partial x} = \lambda(\frac{\partial g}{\partial x})$, $\frac{\partial f}{\partial y} = \lambda(\frac{\partial g}{\partial y})$, etc. Then using combined with the constraint function g to solve λ and the vector that maximizing the original function.

3. the Regularization is essential the LaGrange Multiplier

Stirling Approximation:

$$N! = 1 * 2 * 3 * \dots * N$$

$$\begin{aligned}\log N! &= \log 1 + \log 2 + \log 3 + \dots + \log N \\ &= \sum \log x \\ &= \int \log x \, dx \\ &= x \log x - x \text{ ranged from } N \text{ to } 1 \\ &= N \log N\end{aligned}$$

thus,

$$N! = N^N * e^{-N}$$

Jensen's Inequality:

Given random variable X, $g(X)$ is also random variable, then,

$$E[g(X)] \geq g(E[X])$$

if and only if g is a convex function (**the Hessian $H(g(X)) \geq 0$ for all X**)

More generally, taken latent random variable t, $p(t = a_1) = \beta_1$, $p(t = a_2) = \beta_2$, $p(t = a_3) = \beta_3$, we have

$$\beta_1 \cdot f(a_1) + \beta_2 \cdot f(a_2) + \beta_3 \cdot f(a_3) \geq f(\beta_1 a_1 + \beta_2 a_2 + \beta_3 a_3)$$

where $\beta_1 \cdot f(a_1) + \beta_2 \cdot f(a_2) + \beta_3 \cdot f(a_3)$ is $E_{t \sim p}[f(t)]$, $f(\beta_1 a_1 + \beta_2 a_2 + \beta_3 a_3)$ is $f(E_{t \sim p}[t])$

Mutual Information:

1. **entropy** is the measure of self information $-H(p, p)$.

If we think of a distribution as the tool we use to encode symbols, then entropy measures the number of bits we'll need for this particular distribution.

This is optimal, in that we can't encode the symbols using fewer bits on average.

2. **cross entropy** measures the number of **extra bits** we'll need on average if we encode symbols from y according to $y\hat{}$.

Think of it as a bit tax for encoding symbols from y with an inappropriate distribution $y\hat{}$. It's never negative, and it's 0 only when y and $y\hat{}$ are the same.

minimizing cross entropy = minimizing the KL divergence = minimizing the negative log likelihood of our data.

<http://rdipietro.github.io/friendly-intro-to-cross-entropy-loss/>

3. KL-divergence: the distance of uncertainty between 2 probability distributions

$$\begin{aligned}
 D(P || Q) &= \sum P(x) \log (P(x) / Q(x)) \\
 &= E[P \log (P/Q)] \\
 &= \sum P(x) \log P(x) - \sum P(x) \log Q(x) \\
 &= H(P, P) - H(P, Q) \quad \# \text{ entropy - cross entropy} \\
 &= E[\text{bits of model P}] - E[\text{bits from model P when using model Q}]
 \end{aligned}$$

$$\begin{aligned}
 4. I(X; Y) &= D[p(x, y) || p(x)p(y)] \\
 &= D[p(y|x) || p(y)] \\
 &= D[p(x|y) || p(x)] \\
 &= H(X) + H(Y) - H(X; Y) \\
 &= H(X) - H(X|Y)
 \end{aligned}$$

Thus we can interpret the MI between X and Y as the reduction in

uncertainty about X after observing Y , or, by symmetry, the reduction in uncertainty about Y after observing X.

5. the mutual information can be used to discover interesting relationships between variables in a way that simpler measures, such as correlation coefficient , cannot.

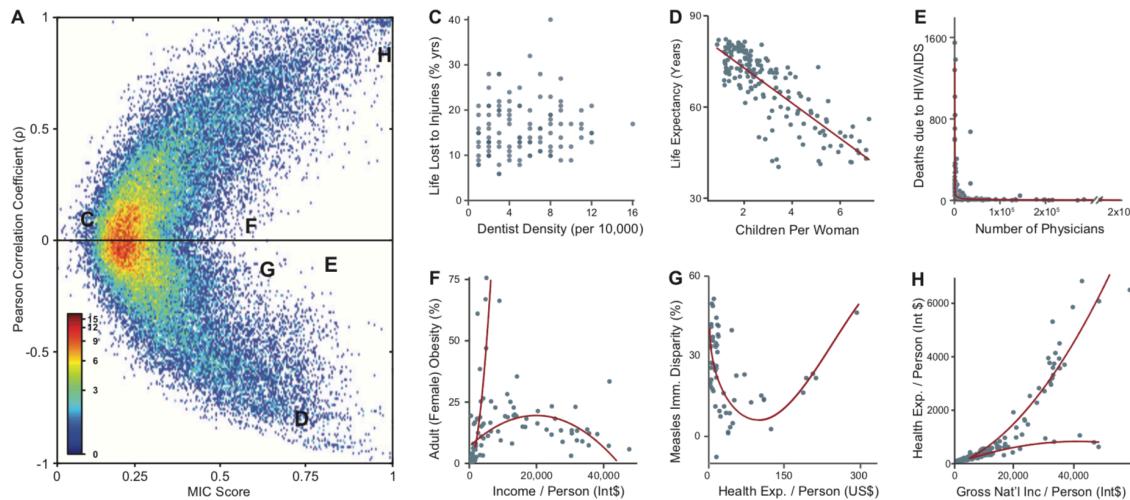


Figure 2.22 Left: Correlation coefficient vs maximal information criterion (MIC) for all pairwise relationships in the WHO data. Right: scatter plots of certain pairs of variables. The red lines are non-parametric smoothing regressions (Section 15.4.6) fit separately to each trend. Source: Figure 4 of (Reshed et al. 2011) . Used with kind permission of David Reshed and the American Association for the Advancement of Science.

Quick notes on Matrix:

$[v_1, v_2]$ is the vertical vector, and $[v_1, v_2]^T$ is the horizontal vector.

Matrix Intuition:

1. Matrix can be either interpreted as a collection of vector or a collection of transformation.
2. Matrix can be also interpreted as the relationship among the high dimensional data. If you dot product 2 vectors XY^T , you will get the covariance matrix of it.

Matrix Determinant:

The determinant of a matrix is the **AREA** of the transformed **Identity** matrix transforming from it.

$$A = [[a \ b]^\top [c \ d]^\top] \quad \det|A| = ad - bc = \prod \lambda_i$$

Matrix Inversion:

the inverted transformation respect to M,

$$M^{-1}M = MM^{-1} = I$$

Matrix Transpose:

$$U^T V = V^T U$$

Singular Matrices:

A Matrix is singular if A^{-1} is undefined.

$$A = [[a \ b]^\top [c \ d]^\top] \quad A^{-1} = 1 / \det|A| [[d \ -b]^\top [-c \ a]^\top]$$

A^{-1} is undefined IFF. the determinant of A ($\det|A|$) is 0.

It can be explained, in another way, as the identity vectors within the Identity matrix become parallel after the transformation.

Matrix Multiplication:

Matrix Multiplication is just the vector transformation.

$$M1 M2 = [[a1, a2][b1, b2]] [[c1, c2][d1, d2]]$$

Think of the 2 X 2 matrix in 2D space with basic vector $I = [i, j] = [[1, 0]$

$[0, 1]$.

In Matrix M1, the vector $[a_1, a_2]$ represent scaling a_1 and a_2 represent to i and j, thus

$$a_1 * [1, 0] + a_2 * [0, 1] = [a_1, a_2]$$

Similarly, the vector $[b_1, b_2]$ represent scaling b_1 and b_2 represent to i and j, thus we have the Matrix M1.

Matrix can be seen as the transformation function, or the collection of transformed basic vector I.

Multiplied by M2, the transformed $[i, j]$ would be scaled by another 2 vectors $[c_1, c_2]$ and $[d_1, d_2]$.

Matrix Dot Product:

$$[a_1, a_2]^T \cdot [b_1, b_2] = a_1 b_1 + a_2 b_2 = |a| |b| \cos\theta$$

The dot product project the 2D space vector $[a_1, a_2]$ on to a 1D line that the $[b_1, b_2]$ represented, and measure the direction (angle) through +/- of its result. $[a_1, a_2]^T \rightarrow [a_1, a_2]$ is called the duality.

If the matrix with column vectors dot product equals 0 to each other, orthogonal to each other in other words, is called the **Orthogonal Matrix**. Furthermore, if the length column vectors are 1, then the matrix is called **Orthonormal Matrix**

Vector Projection

$$\text{Proj}_v(u) = [(u^T v) / \|v\|^2] v$$

as projecting vector u onto vector y .

Eigenvector and Eigenvalue:

Any Symmetric Matrix ($M = M^T$) can be written as

$$\mathbf{M} = \mathbf{V}\Lambda\mathbf{V}^{-1} = \sum \lambda_i \mathbf{v}\mathbf{v}^{-1}$$

where

\mathbf{V} is a collection of (orthogonal or so-called rotation) eigenvectors respect to eigenvalues Λ , $\mathbf{V}^{-1} = \mathbf{V}^T$

Λ is the diagonal matrix of eigenvalues λ_i , which is

$$\mathbf{M}\mathbf{V} = \Lambda\mathbf{V}$$

<https://www.youtube.com/watch?v=spNpfmWZBmg>

Affine (Linear) Transformation:

If the sum of the transformation equal to the transformation of the sum, then it is called the Affine (Linear) Transformation.

<https://www.youtube.com/watch?v=y2WqIXrjyC4&index=2&list=PLIXfTHzgMRUIqYrutsFXCOmiqKUgOgGJ5>

Singular Value Decomposition

For any matrix $A \in \mathbb{R}^{m \times n}$ there exist orthogonal matrices $U \in \mathbb{R}^{m \times m}$, $V \in \mathbb{R}^{n \times n}$ and a 'diagonal' matrix $\Sigma \in \mathbb{R}^{m \times n}$, i.e.,

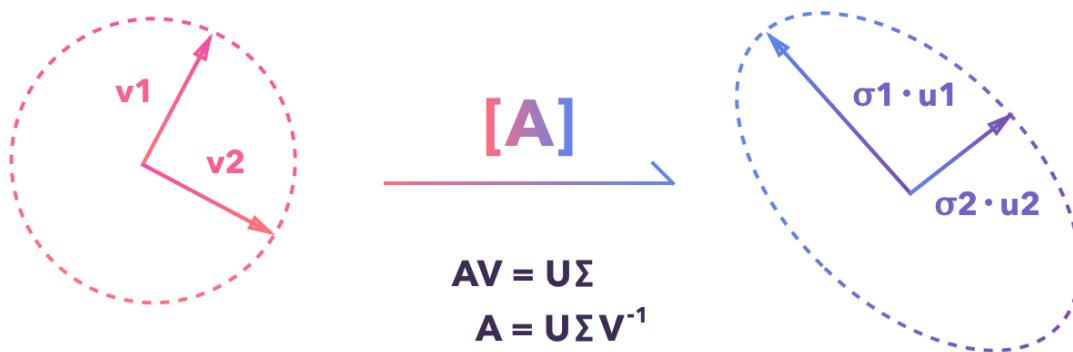
$$\Sigma = \begin{pmatrix} \sigma_1 & & & \\ & \ddots & & \\ & & \sigma_r & \\ & & & 0 \\ & & & & \ddots \\ 0 & & & & & 0 \\ \vdots & & & & & \vdots \\ 0 & & & & & 0 \end{pmatrix} \quad \text{for } m \geq n$$

with diagonal entries

$$\sigma_1 \geq \dots \geq \sigma_r > \sigma_{r+1} = \dots = \sigma_{\min\{m,n\}} = 0$$

such that $A = U\Sigma V^T$

Think of a matrix A that transforming from a orthogonal matrix $V = [v_1, v_2]$ to another orthogonal matrix $U = [u_1, u_2]$ as of transforming from a circle with radius v_1 ($v_1 = v_2$ in this special case) to a ellipse with axis u_1, u_2 .



where

- **U** is the eigenvectors respect to the squared root of eigenvalues Σ for covariance matrix AA^T
- Say A is the data (design) matrix, the k -th principle components of A are $u_1\sigma_1v_1^T, u_2\sigma_2v_2^T, \dots, u_k\sigma_kv_k^T$
- To projecting N demential data matrix A onto k demential principle components, $A_k = U^TA [1:k]$

- Say A is the data (design) matrix, recall that the unbiased estimator of Variance is $\text{Var}(A) = (1 / (N - 1)) \text{AA}^T$, then $\text{Var}(A) = (1 / (N - 1)) \text{U}\Sigma^2\text{U}^T$

Principle Component Analysis and Mahalanobis Distance:

Principle Component Analysis:

PCA reduces the demotions of the data while remaining the most important information.

It attempt to disentangle the unknown factors of variation underlying the data , this disentangling takes the form of finding a rotation of the input space (described by U) that aligns the principal axes of variance with the basis of the new representation space associated with $\mathbf{z} = \mathbf{x}^T \mathbf{U}$

To perform PCA, you calculate the eigenvectors of the data's covariance matrix. Typically, we use SVD to decompose the data then retrieving the **singular vectors U and variance Σ** as the eigenvectors and squared root of eigenvalues of the covariance matrix.

Mahalanobis Distance:

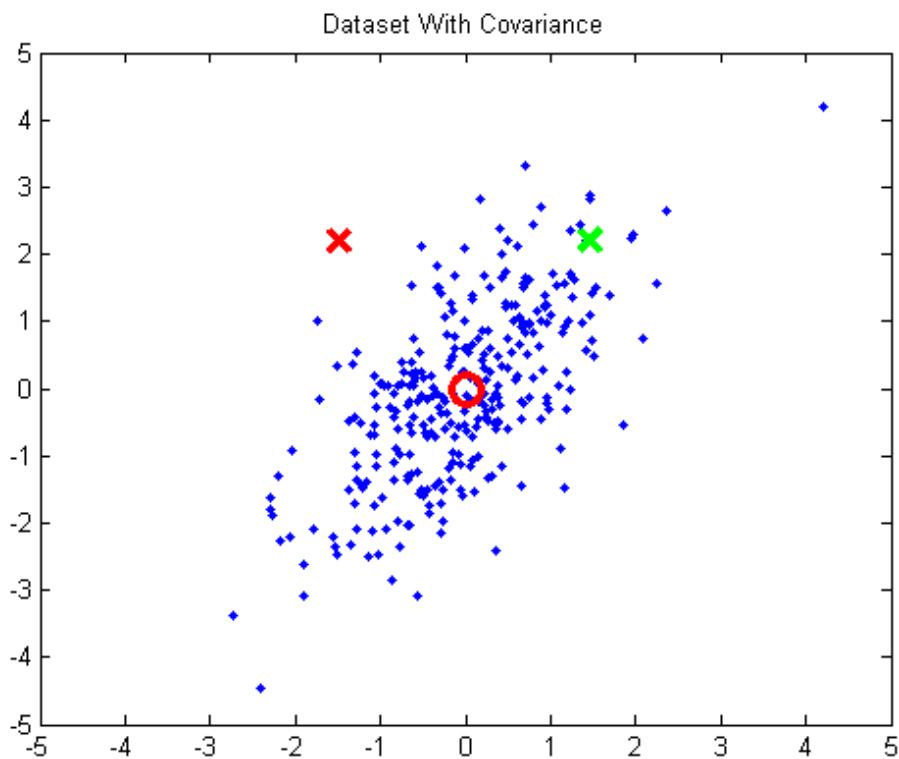
$$d_M(x, y) = \sqrt{(x - y)^T S^{-1} (x - y)}$$

where S^{-1} is the inversion of the covariance matrix of X and Y.

[Side note]: the probability density function for a multivariate Gaussian distribution uses the Mahalanobis distance instead of the Euclidean.

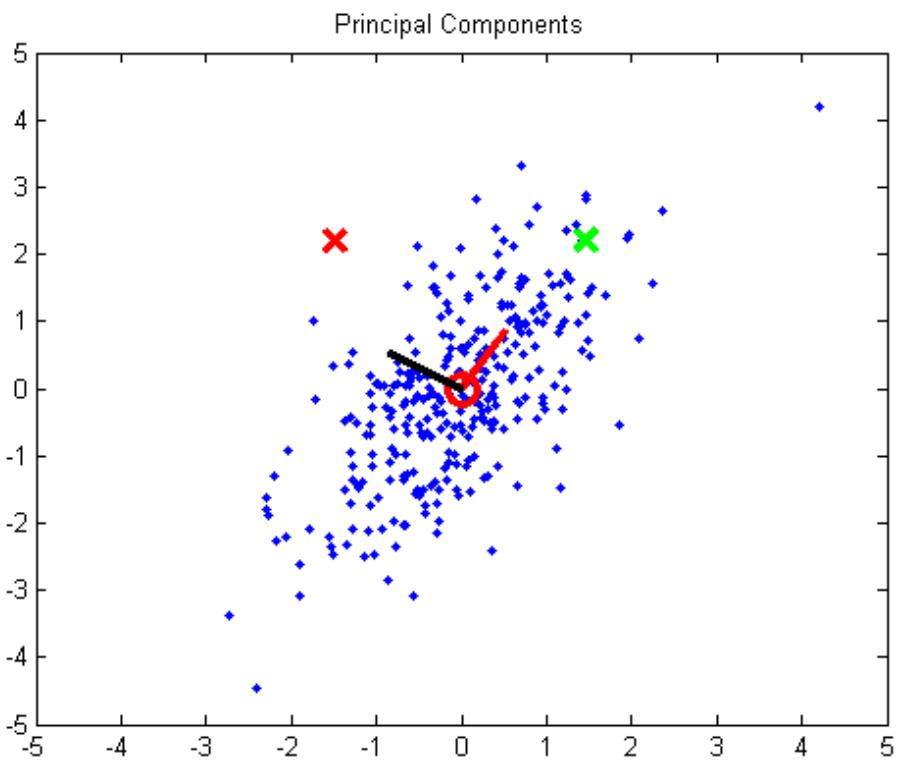
Mahalanobis Distance measuring the Euclidean Distance throughout the Non-Correlational data (wipe out the correlation through S^{-1}).

Intuitively, one can think of it as the following steps:

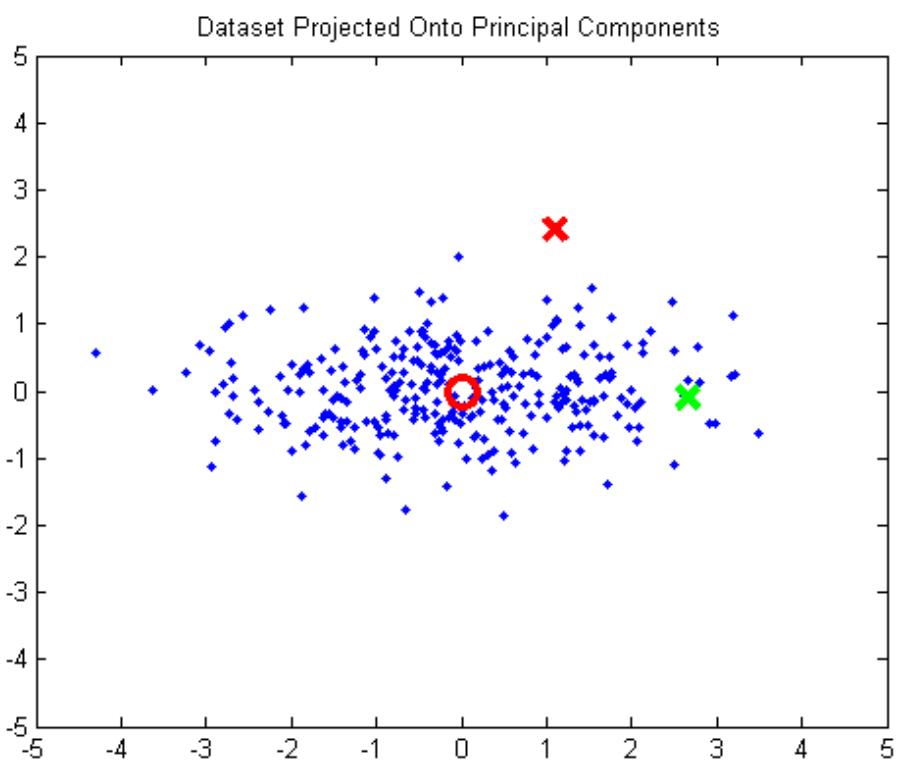


Here, the covariance of the data is non-zero since the principle components are **not aligned with axes**. Thus the RED and GREEN test points shouldn't be considered equally with the same similarity to the data distribution even though they have the same Euclidean Distance to the mean of the data.

The solution is to find the principle components of the data using SVD:

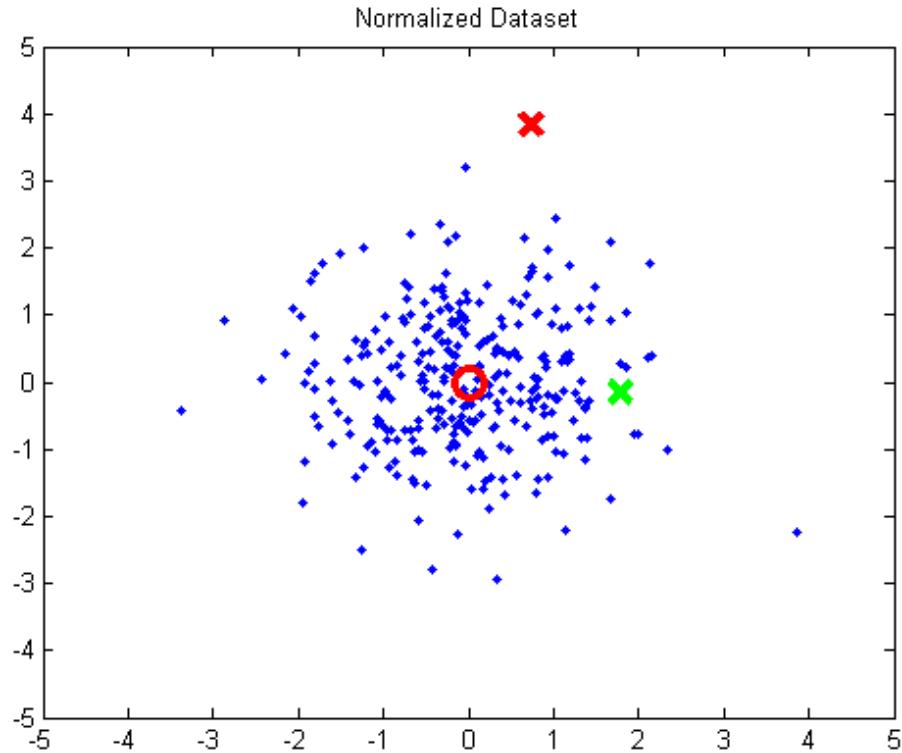


Then projecting data onto these principle components. This is the crucial step to wipe out the covariance of the data.



Finally, normalizing the data along with each of the principle components by

dividing the standard deviation Σ (the squared root of eigenvalues for the covariance matrix). The Euclidean Distance of GREEE point is lower than the RED one now.



The 3 steps above are exactly the **PCA whitening**.

$$x_{\text{PCAwhite},i} = \frac{x_{\text{rot},i}}{\sqrt{\lambda_i + \epsilon}}.$$

Although calculating the Mahalanobis distance between our two example points yields a different value than calculating the Euclidean distance between the PCA Whitened example points, it's still the intuitive way of understanding how correlation needs to be taken into account for distances.

<http://mccormickml.com/2014/07/21/mahalanobis-distance/>

What is the probability of data in machine learning exactly ?

assuming all possible values of an input, i.e. the price of the house, fits

into some probability distribution parameterized by θ , which is the probability of the data $p(x | \theta)$. for a pixel from an image, it is always chosen from [0, 255] with some probability distribution underneath.

Generalization Error:

1. the error (**expected loss**) for the test data is the probability distribution that the hypothesis $h\theta$ is unequal to the true outcome y

$$\varepsilon(h) = P \sim D(h\theta \neq y)$$

2. minimizing the generalization error is equivalent to maximizing the likelihood.

Bias-Variance tradeoff

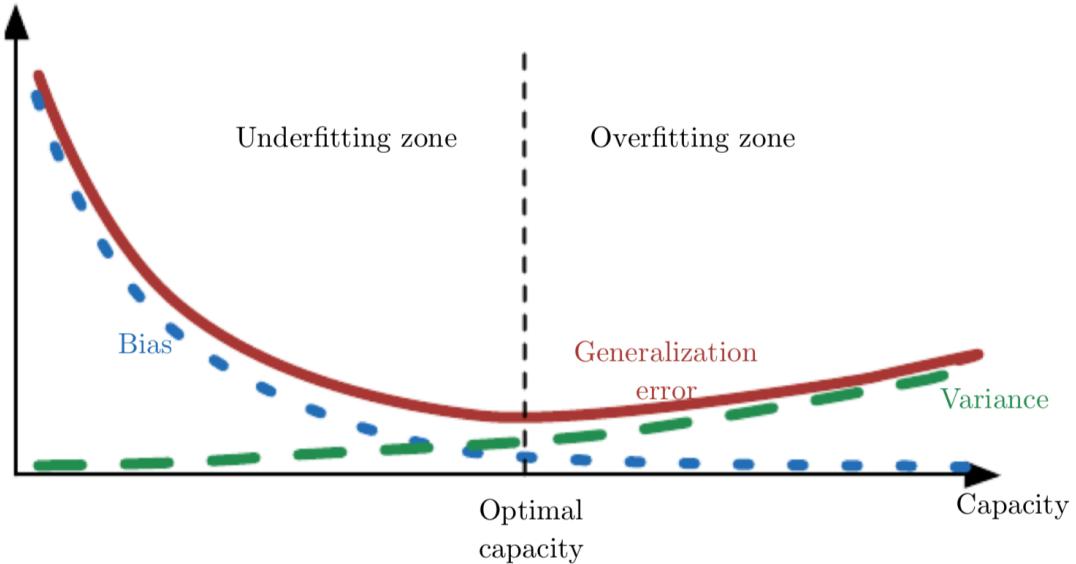
for vary range of loss function, taking the expectation over the data set, we have

$$\begin{aligned}
 & \text{expected loss (Mean Square Error for example)} \\
 &= E[(\theta_m - \theta)^2] \\
 &= E[\theta_m^2 - 2\theta \cdot \theta_m + \theta^2] \\
 &= E[\theta_m^2] - 2\theta \cdot E[\theta_m] + \theta^2 \# E[\theta^2] = \theta^2 \\
 &= (E[\theta_m]^2 - 2\theta \cdot E[\theta_m] + \theta^2) + (E[\theta_m^2] - E[\theta_m]^2) \\
 &= (E[\theta_m] - \theta)^2 + (E[\theta_m^2] - E[\theta_m]^2)
 \end{aligned}$$

where θ_m is the parameter of the sample distribution.

By the definition, Bias is just the distance between the expected empirical distribution and the true distribution, thus,

$$MSE = \text{Bias}(\theta_m)^2 + \text{Var}(\theta_m)$$



In the case where generalization error is measured by the MSE (where bias and variance are meaningful components of generalization error), **increasing capacity tends to increase variance and decrease bias.**

see also **Deep Learning Book 5.4.4**

from the perspective of learning theory, we have

$$\epsilon(\text{hypotheses-hat}) \leq (\min \epsilon(\text{hypotheses})) + 2((1/2m) \log(2k/\gamma))$$

where

ϵ is the generalization error,

m is the training size,

k is the number of hypotheses,

hypotheses-hat is the hypotheses of new example drawing from the training set,

$\min \epsilon(\text{hypotheses})$ is the bias,

$2((1/2m) \log(2k/\delta))$ is the variance.

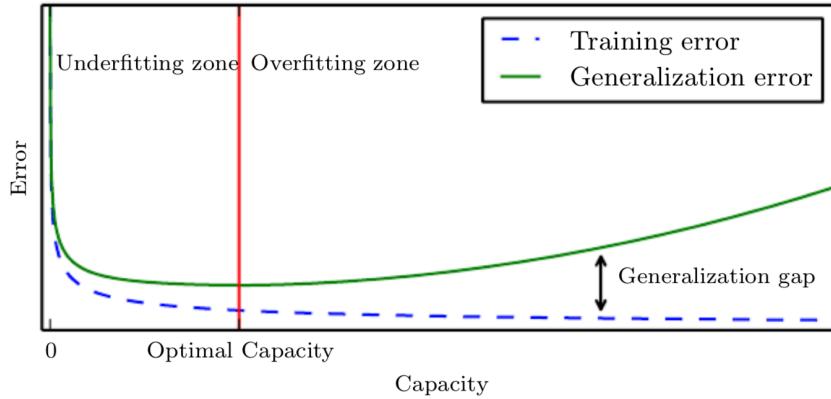


Figure 5.3: Typical relationship between capacity and error. Training and test error behave differently. At the left end of the graph, training error and generalization error are both high. This is the *underfitting regime*. As we increase capacity, training error decreases, but the gap between training and generalization error increases. Eventually, the size of this gap outweighs the decrease in training error, and we enter the *overfitting regime*, where capacity is too large, above the *optimal capacity*.

Big O Notation:

1. Definition: How **TIME** scales respect to **INPUT**.
For example, a for-loop takes $O(N)$ amount of time, a for-loop inside the identical loop takes $O(N^2)$ amount of time, etc.
2. Constant doesn't matter
 $O(2N) \approx O(N)$
as N goes infinity.
3. Dropping the non-dominant term
 $O(N) + O(N^2) = O(N + N^2) \approx O(N)$
since $O(N)$ is non-dominant compared to $O(N^2)$

Monte Carlo Approximation:

The basic idea of Monte Carlo Approximation is that, for complex probability distribution, it's unable to compute the expected value. Instead we can use Monte Carlo Approximation to approximate the expected value using the average (mean) value of samples that drawing from the distribution.

$$\mu = (1/N) \sum f(x_i), \quad x \text{ ranging from } 1 \dots N$$

For the Bayesian Inferences, it's usually unable to compute the marginal likelihood which is a NP hard problem, so instead another analyzable probability distribution could be introduced to address the problem.

Say we have the marginal likelihood Z , and we introduce a gaussian distribution $Q(\theta)$

$$Z = \int P(y|\theta) P(\theta) d\theta \\ = \int P(y|\theta) P(\theta) Q(\theta) / Q(\theta) d\theta$$

Rewritten the ratio of the likelihood against the introduced distribution as $W(\theta)$, we have

$$Z = \int W(\theta) Q(\theta) d\theta$$

then the marginal likelihood is converted to the Expectation of $W(\theta)$ with the probability of $Q(\theta)$ over θ based on **the LOTUS law**. Thus with the **Law of Large Number**, we have

$$Z = (1/N) \sum_i W(\theta^i)$$

as i goes infinity.

Markov Chain:

0. Taken P as the Markov Chain(Matrix)

1. Markov Property

the future is only dependent on the current.

$$P(T_{n+1} | T_n, T_{n-1}, T_{n-2}, \dots, T_1) = P(T_{n+1} | T_n)$$

$$2. P(Tm + Tn) = P(Tm) P(Tn)$$

$$\begin{aligned}
 3. \frac{dP(T)}{dT} &= \lim_{dT \rightarrow 0} (P(T + dT) - P(T)) / dT \\
 &= \lim_{dT \rightarrow 0} (P(T)P(dT) - P(T)) / dt \quad \# P(Tm + Tn) = \\
 P(Tm)P(Tn) &= P(T) \lim_{dT \rightarrow 0} (P(dT) - I) / dT \quad \# I \text{ is the identity} \\
 \text{matrix} &= P(T) Q(dt)
 \end{aligned}$$

As the derivative $dP(T) / dT$ going to be 0, **P(T)** would not to be **changed with the Time**, which means **P(T)** will going to be the stationary distribution π .

4. Ergodic

A Markov chain is ergodic if it has limiting stationary distribution π and if none of the elements of this distribution are zero.

Softmax function and Temperature:

Softmax function is the generalized function of Logistic Regression. It can be interpreted as the posterior probability $p(y = c | x ; \theta)$ that inferred by the likelihood of data x modeled by the exponential family distribution (1), Gaussian, Bernoulli, etc., divided by the marginal likelihood of (1), that is

$$p(y = c | x ; \theta) = S(\eta_i) = e^{\eta_i} / \sum e^{\eta}$$

where η is the nature parameter $\eta_i = k(\theta^T x)$.

In terms of thermodynamics statistical, $-\eta$ is called the **Energy function**, and $\sum e^{\eta}$ is called the **partition function**. There we can introduce another parameter T which represents the **Temperature** of the system.

$$S(\eta_i/T) = e^{\eta_i/T} / \sum e^{\eta_j/T}$$

As T goes to infinity, every $e^{\eta_i/T}$ goes to 1, then $S(\eta_i/T) \approx 1/(1+1+\dots)$ becomes a uniform distribution which owns the highest entropy, which implies the higher the Temperature, the more uncertainty the system.

Otherwise as T goes to 0, $e^{\eta_i/T}$ goes to infinity, then $S(\eta_i/T) \approx 1$, and every other $e^{\eta_j/T}$ goes to 0 since the denominator should sum to 1, which implies the lower the Temperature, the more certainty the system.

The form of $e^{\beta\eta_i} / \sum e^{\beta\eta_j}$ is called the **Boltzmann Distribution**.

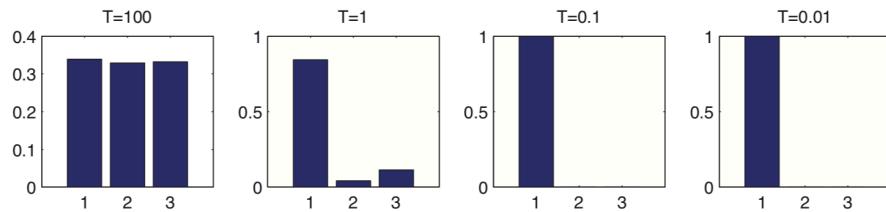


Figure 4.4 Softmax distribution $S(\eta/T)$, where $\eta = (3, 0, 1)$, at different temperatures T . When the temperature is high (left), the distribution is uniform, whereas when the temperature is low (right), the distribution is “spiky”, with all its mass on the largest element. Figure generated by softmaxDemo2.

Support Vector Machine:

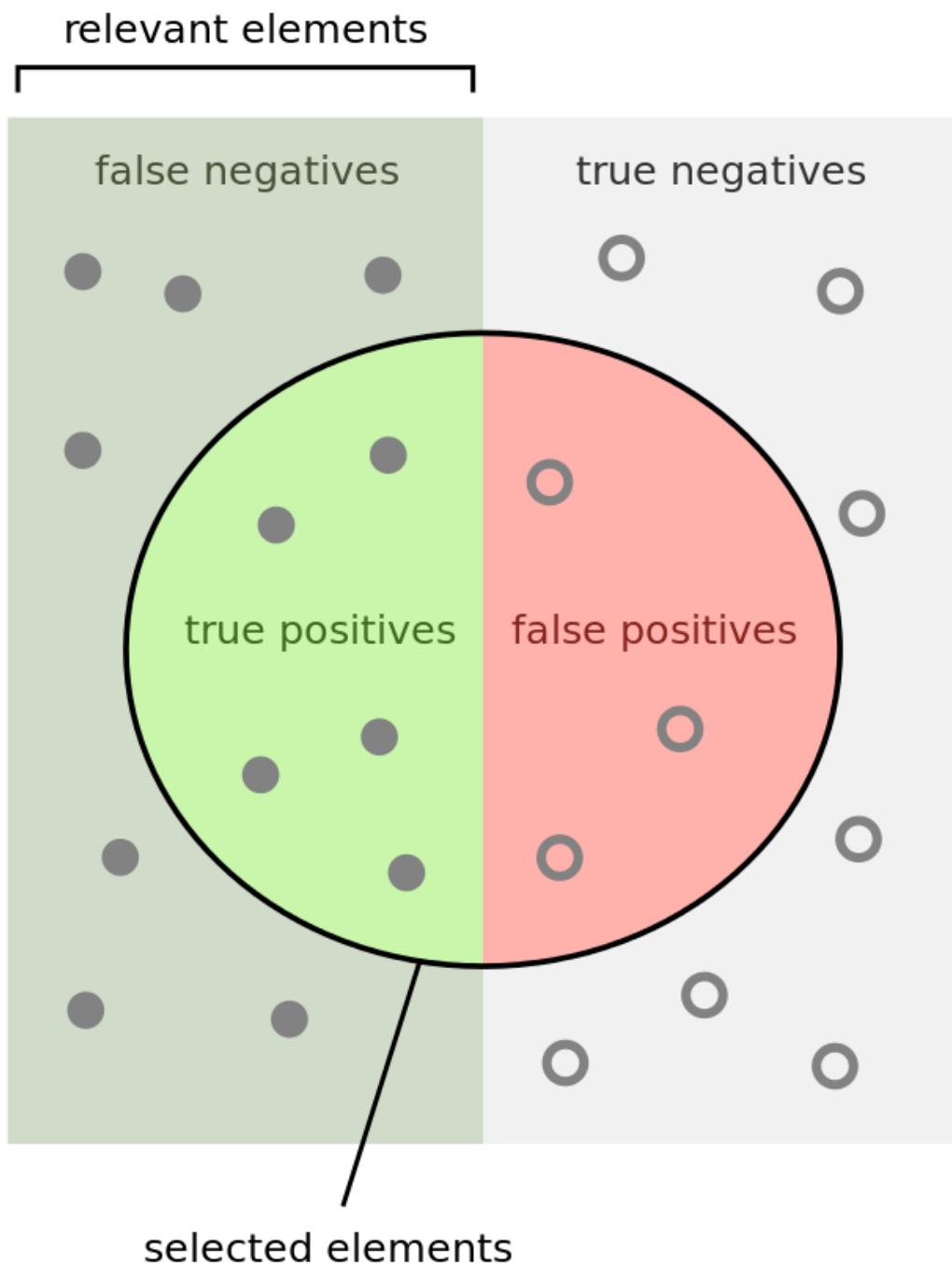
Precision / Recall:

Suppose a computer program for recognizing dogs in photographs identifies eight dogs in a picture containing 12 dogs and some cats. Of the eight dogs identified, five actually are dogs (true positives), while the rest are cats (false positives). The program's precision is 5/8 while its recall is 5/12.

When a search engine returns 30 pages only 20 of which were relevant while failing to return 40 additional relevant pages, its precision is $20/30 = 2/3$ while its recall is $20/60 = 1/3$.

Precision is "how useful the search results are", and recall is "how complete the results are"

If you want to increase the recall, the precision decreases.



How many selected items are relevant?

$$\text{Precision} = \frac{\text{Number of relevant items}}{\text{Number of selected items}}$$

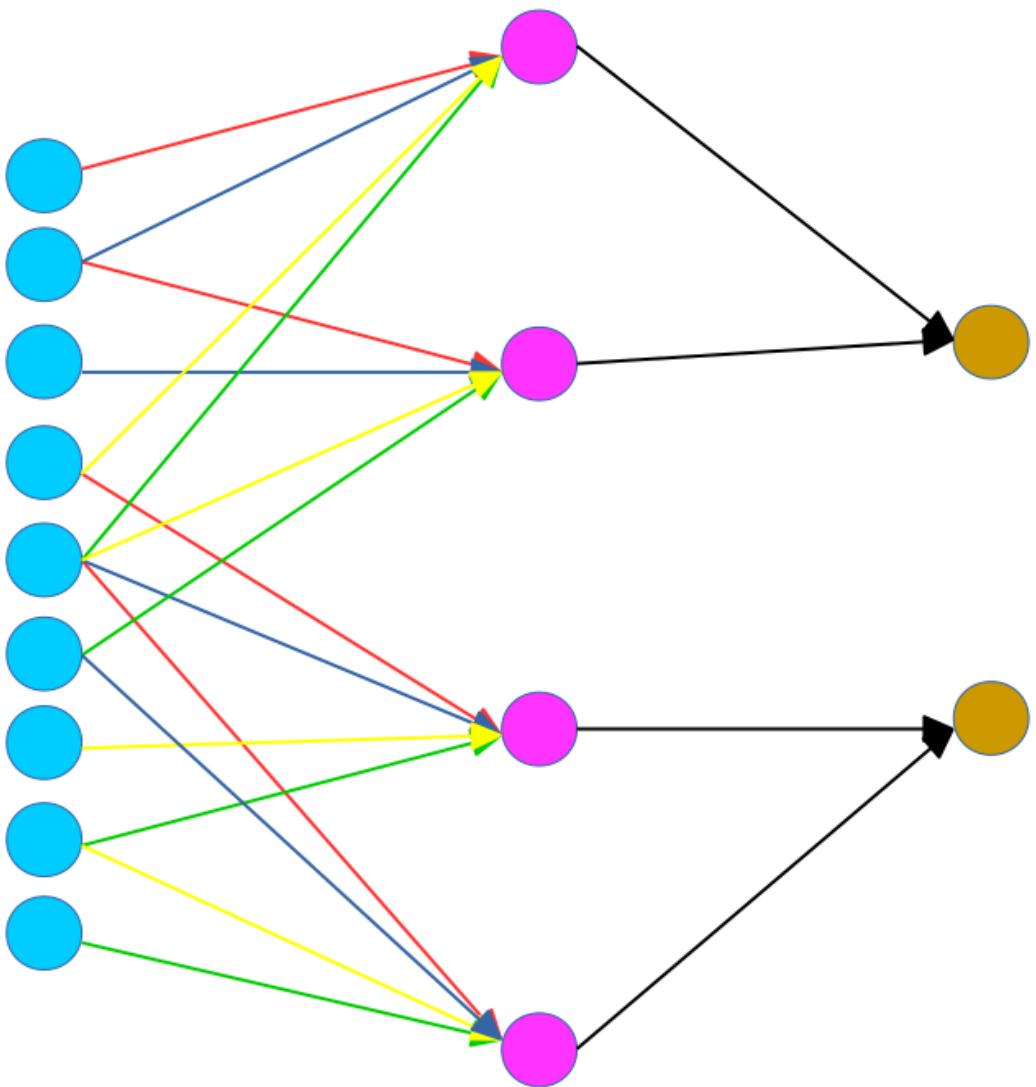

How many relevant items are selected?

$$\text{Recall} = \frac{\text{Number of relevant items}}{\text{Number of all relevant items}}$$


What is Precision-Recall (PR) curve?

<https://www.quora.com/What-is-Precision-Recall-PR-curve>

CNN:



$$\begin{array}{c}
 \text{Input Feature Map} \\
 \begin{matrix} \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} \end{matrix} \\
 \times \quad \begin{matrix} \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} \end{matrix} = \quad \begin{matrix} \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} \end{matrix} \quad | \quad \begin{matrix} \text{---} & \text{---} & \text{---} \\ \text{---} & \text{---} & \text{---} \end{matrix} \\
 \text{Kernel} \qquad \qquad \qquad \text{Conv.} \qquad \qquad \qquad \text{Pool}
 \end{array}$$

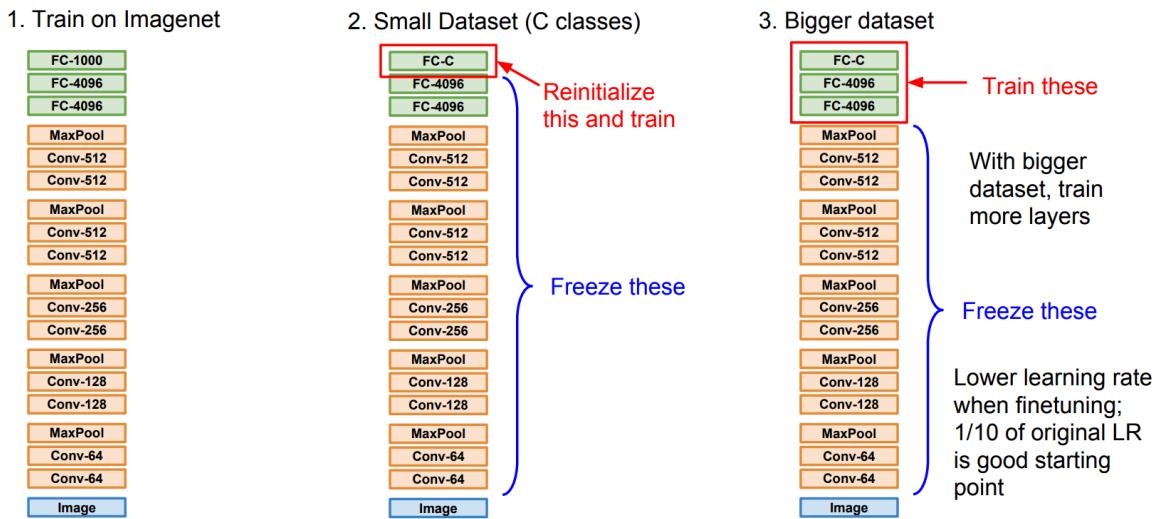
When calculating the derivative of \mathbf{W}_{red} , the **4 red lines** should be taken into account in contrast with **1 line** for normal DNN.

After the max pooling layer, some of the magenta neurons would be dismissed as of the Regularization or Dropout for downsampling or preventing overfitting, say the top and bottom neurons, thus the weights directed to those neurons would be dismissed as well.

for more details, [Backpropagation In Convolutional Neural Networks](#)

Transfer Learning:

Transfer Learning with CNNs



Questions:

1. What is the **Universal Approximation Theorem**? A: <http://neuralnetworksanddeeplearning.com/chap4.html>
2. what the relationship between Naive Bayesian and Logistic Regression?
3. parametric vs non-parametric. parametric: like linear regression with weight parameters w; non-parametric: like nearest neighbor regression.
4. We run it through L2- normalization, then an affine transformation step, then another normalization step for mapping the result to a normalized feature space. For the affine transform we consider a learned fully-connected layer. **[Paper] Scalable Object Detection for Stylized Objects**
5. The PCA projection can be implemented with a shifting and a fully connected (FC) layer **[Paper] Deep Image Retrieval - Learning global representations for image search**
6. **DATA - VECTOR - Probability**
7. relationship between svm and svd

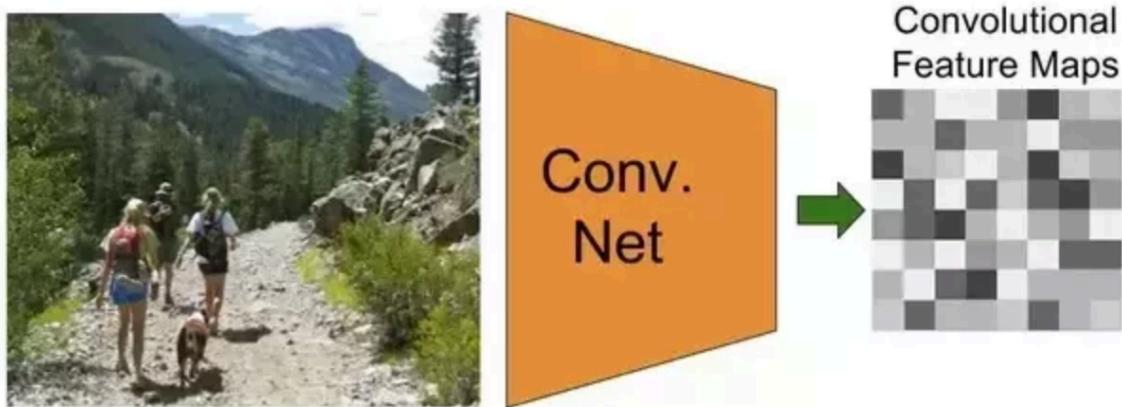
Paper Notes:

1. [Leveraging category-level labels for instance-level image retrieval] finetuning the pretrained models for retrieval can bring a significant improvement
2. [Learning fine-grained image similarity with deep ranking] One can distinguish offline hard negative mining, where a second model is trained to identify hard cases on the whole dataset.
3. [Facenet: A unified embedding for face recognition and clustering] suggest various methods for picking the negative samples.

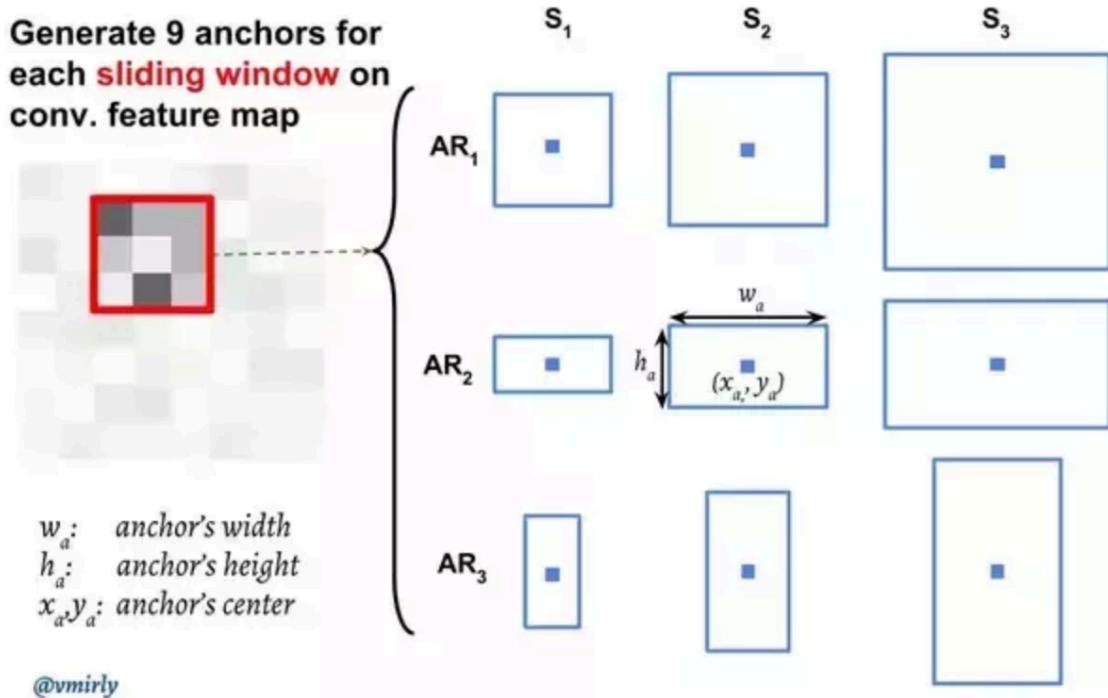
Miscellaneous:

1. ROI Pooling
2. How does the region proposal network (RPN) in Faster R-CNN work?

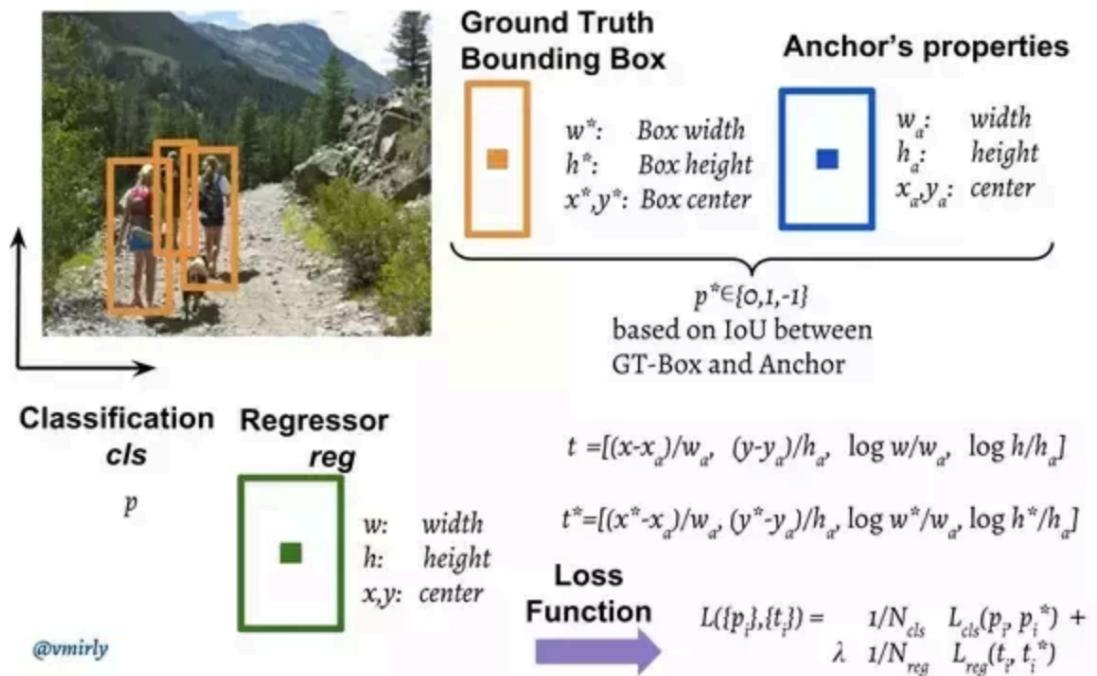
1. In the first step, the input image goes through a convolution network which will output a set of convolutional feature maps on the last convolutional layer:



2. Then a sliding window is run spatially on these feature maps. The size of sliding window is $n \times n$ (here 3×3). For each sliding window, a set of 9 anchors are generated which all have the same center (x_a, y_a) but with 3 different aspect ratios and 3 different scales as shown below. Note that all these coordinates are computed with respect to the original image.



3. Finally, the 3×3 spatial features extracted from those convolution feature maps (shown above within **red box**) are fed to a smaller network which has two tasks: classification (cls) and regression (reg). The output of regressor determines a predicted bounding-box (x, y, w, h) . The output of classification sub-network is a probability p indicating whether the predicted box contains an object (1) or it is from background (0 for no object).



3. **How to train Fast-RCNN ?** First RPN is trained and regions obtained are used to train detection network. Then RPN is retrained but this time with the convolution part initialized with weights from convolution part of detection network.
4. Different between **latent variables** of probabilistic models and **hidden unit** of neural network. <https://www.quora.com/How-are-latent-factor-neural-networks-trained>