

Vector Algebra
 $XX^T = \sum_{i=1}^n x_i x_i^T = \sum_{i=1}^n x_i^2 = \|X\|_F^2$
 $\frac{\partial}{\partial x} (x^T x) = \frac{\partial}{\partial x} (x^T x) = 2x$
 $\frac{\partial}{\partial x} (x^T A x) = (A + A^T)x$
 Remember: If hard derivative, element-wise.
 Trick: derive wrt. every vector element, many summands become 0.

Convex
 A twice differentiable convex function is its Hessian H_{xx} satisfies $H_{xx} \succeq 0$.
 The function would be concave if $-f(x)$ is strictly convex.
 Alternative criterion: convex if $\forall x, y, \lambda \in [0, 1]$
 $f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y)$
 i.e. if the function lies below the connecting line at x, y .
 convex + concave = constant

Random Formula Knowledge
 $\int_{-\infty}^{\infty} \exp\left\{-\frac{1}{2}(x-\mu)^2\right\} dx = \sqrt{2\pi}$
 $\int_{-\infty}^{\infty} \exp\left\{-a\left[\left(\frac{x}{b}\right)^2 - \frac{c}{a}\right]\right\} dx = \frac{\sqrt{2\pi}}{b} \exp\left\{\frac{c}{a}\right\}$
 $\log(a/b) = \ln(a/b)$
 $\log(a^b) = b \cdot \log(a)$

Covariance Matrix, Bivariate normal random variables
 $N(\Sigma, \mu) \Rightarrow f(x) = \frac{1}{(2\pi)^{n/2} |\Sigma|} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right\}$
 $\Sigma = \text{Cov}(X) = E[(X-\mu)(X-\mu)^T]$
 $= \begin{pmatrix} \text{Var}(x_1) & \dots & \text{Cov}(x_1, x_n) \\ \vdots & \ddots & \vdots \\ \text{Cov}(x_n, x_1) & \dots & \text{Var}(x_n) \end{pmatrix}$

PCA & Empirical Covariance Matrix
 Center Data by subtracting the mean. Compute $\Sigma = \frac{1}{n} \sum_{i=1}^n x_i x_i^T$ and Eigenvalues $\det(\Sigma - \lambda I) = 0$ and eigenvectors $(\Sigma - \lambda I)x = 0$.
 For solving: Gaussian Elimination (same change) Determinants are computed by taking $a_{ij} \cdot \det(\text{Matrix without row } i \text{ and col } j)$
 $-a_{11} a_{22} \dots + a_{12} a_{21} \dots$ etc.
 And then make the result unit vectors for PCA. Inverse of a matrix

Write down M^{-1} . Apply row-wise operators to both sides until LHS is I .
Confusion Matrix $\text{rec} = \frac{TP}{TP+FP}$
 Table with numbers of actual vs predictions
 Recall = $P(\text{is Dog} | \text{pred Dog})$. Precision = $P(\text{pred Dog} | \text{is Dog})$
 $F1 = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$

Linear Regression and Regularizers
 $\arg \min_w \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|_1$ (Lasso)
 or $\sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|_2$ (Ridge)
 L2 drives params down quickly but doesn't zero them out.
 Elastic net

Gradient Descent
 $x_{\min} = x_n - \gamma \nabla f(x_n)$, step size $\gamma_n = \frac{1}{n}$
 $\gamma_n = \frac{(x_n - x_{n-1})^T (\nabla f(x_n) - \nabla f(x_{n-1}))}{\|\nabla f(x_n) - \nabla f(x_{n-1})\|^2}$

Logistic Regression
 $\ln\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$
 $\log = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n$
 odds that $Y=1 = e^{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n}$

Generative & Discriminative Models
 Gen: $P(x, y)$ and $P(y)$ model how data was created $\Rightarrow P(x|y)$ example: Naive Bayes
 Dis: $P(y|x)$ models decision boundary between classes
 e.g. SVM, K-Means Clustering, Perceptron

ROC Curve
 True Positives (as ratio of all positives on x this is sensitivity) how sensitive is the test to sick people?
 True Negatives ... on x this is specificity. Increase in one, decrease in the other.
 Test (Classifier) determines those depending on some threshold that can be moved.
 ROC Curve is sensitivity on x -axis and (1-Specificity) on y -axis, i.e. the false positives. Better classifiers have \uparrow and not \downarrow .

Naive Bayes
 Conditionally assumes independence of feature dimensions given y :
 $P(x_1, y) = P(x_1|y)P(x_2|y)P(x_3|y) \dots$
 and compute $P(x|X) = P(x_1|x)P(x_2|x)P(x_3|x) \dots$
 Maximum A Posteriori (MAP) $P(x|X)P(x) / P(x)$
 $P(\theta|x)$ maximized $\Rightarrow \max P(x|\theta)P(\theta)$
 Maximum Likelihood Estimation (MLE)
 $\max P(x|\theta)$
 Perceptron Algorithm

Gradient descent on perceptron loss function.
 $\nabla_{\theta} (w, x, y) = \begin{cases} 0 & \text{if } y = \text{sign}(w^T x) \\ -x & \text{if } y \neq \text{sign}(w^T x) \end{cases}$
 Logistic smoother: Counting words... over is not a defect $\Rightarrow \text{prob} = 0 \Rightarrow$ problem.
 Solution: Add 1 to every word's occurrence count in the data, and add 1 through the divider. $\frac{1}{n+1}$

Perceptron and SVM
 AA is always Kernel Requirements: Mat: $K(x, x')$ is also a Gram matrix.
 Pos. semi-def. $\rightarrow K(x, x') = k(x, x')$ Gram Matrix K must be pos. semi-definite. $k(x_0, x_0) \geq 0$.
 $k(x_0, x_0) = \sum_{i=1}^n x_{0i}^2 = \|x_0\|_2^2 \geq 0$
 Most also be symmetric and pos. semi-definite.
 So $k_{12} = k_{21}$ is not necessarily a kernel. $k_{12} = k_{21}$ is not necessarily a kernel. $k_{12} = k_{21}$ is not necessarily a kernel.
 For det: $K \succeq 0$. EM Algorithm
 Have Data points but don't know their color or distribution. Assume Mixture Distribution. To guess some parameters. 2. Expect those PDFs of every data point to get likelihood of red/blue for data point given the guessed parameters.
 3. create weights $w = \frac{\text{red}}{\text{red} + \text{blue}}$ and send for blue.
 4. Compute mean and variance for red and blue PDF. \Rightarrow New params for red and blue PDF.
 estimate = $\text{sum}(\text{weights} * (\text{data} - \text{old Mean}))$
 estimateMean = $\text{sum}(\text{data} * \text{weights}) / \text{sum}(\text{weights})$
 Neural Network
 Given $\phi(z) = \frac{1}{1 + \exp(-z)}$ value, layer $\Rightarrow \phi(w_0 + \sum w_j x_j)$
 Treat final output as 1 if greater than 0.5, otherwise 0.
 Give weights for each output is not necessarily 0.5, otherwise 0.
 OR: $A \cdot B = \text{round}(\phi(w_0 + w_1 x_1 + w_2 x_2)) \Rightarrow w = [0.5, 1]^T$
 $A \cdot B = [0.5, 1]^T$ (for this ϕ)
 Task: Create a MLP with only one hidden layer and an arbitrary number of nodes that implements XORs. Possible approach: turn expr into ORs of ANDs and try patching all ORs into second layer.
 1st layer: $(A \cdot B) \Rightarrow h_1 = \phi(A + B \cdot C - 2.5)$ To get close enough to 1.
 $(A \cdot B) \Rightarrow h_2 = \phi(B \cdot C - D - 1.5)$ To get close enough to 0.
 $(A \cdot B) \Rightarrow h_3 = \phi(-A + B \cdot C - 0.4)$ To get close enough to 1.
 Task: XOR with as few units as possible.
 $\text{out} = \phi(-\phi(A+B) + 0.5 \cdot \phi(2A+2B))$
 $\text{out} = 0.84$

Perceptron: $\min \sum \max\{0, -x^T w\}$
SVM: $\min \sum \max\{0, 1 - x^T w\}$
Gradient $\nabla G(w) = \frac{1}{n} \sum \nabla G_i(x_i, w)$
 $\nabla G_i(x, w) = \nabla \max\{0, 1 - x^T w\} = \nabla \lambda \|w\|_2$
 $(\nabla \lambda \|w\|_2)_j = \frac{2\lambda w_j}{\|w\|_2} = 2\lambda y_j \Rightarrow \lambda \|w\|_2^2 = 1$
 $\nabla \max\{0, 1 - x^T w\} = \begin{cases} -x & \text{if } x^T w < 1 \\ 0 & \text{otherwise} \end{cases}$
 $= -\sum x_i x_i$
 $\therefore \nabla \text{sign}(w^T x_i)$

SVM: To get 0 Loss, you have to classify not only correctly, but correctly by some margin. This is hinge Loss: just shifted perceptron loss by 1. If it were some other nonzero number, we would get the same solution, just by scaling w . So we also want to keep the weights small to keep the margin larger, which scales with $\frac{1}{\|w\|}$.

SVM Gradient Update Derivation with SGD
 From above Gradient, we get $w = \begin{cases} w + \eta(1-2xw) & \text{if } x^T w < 1 \\ w + \eta(x - x^T w) & \text{otherwise} \end{cases}$
 Stochastic Gradient Descent (SGD) with backsize
 Initialize somehow, then: for $t = 1, 2, \dots, \text{dag}$
 pick $i \sim \text{Unif}(1, n)$.
 if $\frac{1}{2} \text{sign}(w^T x_i) \neq y_i$ {
 $w_{\text{bin}} = w + \eta y_i x_i$
 $\text{else } w_{\text{bin}} = w - \eta y_i x_i$

Kernel Trick
 wants: artificially high dimensionality to linearly separate data that isn't linearly separable.
 $\phi(x) = \begin{pmatrix} x_1^2 \\ x_1 x_2 \\ x_2^2 \end{pmatrix} \Rightarrow \phi(x)^T \phi(x') = (x_1^2 x_1'^2 + 2x_1 x_1' x_2 x_2' + x_2^2 x_2'^2) = (x_1 x_1' + x_2 x_2')^2$
 It is easier to compute than first transforming x_1 and x_2 with $\phi(x)$. Generalization with d input dimensions: input $\phi(x) = [x_1, x_2, \dots, x_n]^T$ same for x' , both are inputs. But $(x^T x')^2 \in \Theta(d^2)$ is way too costly to compute.
 For monomials of degree m instead of 2: $\phi(x) = [x_1^m, x_1^{m-1} x_2, \dots, x_1 x_2^{m-1}, x_2^m]^T$
 All monomials up to degree m : $k(x, x') = (1 + x^T x')^m$ would be exponential in m .
 Explicit computation would be exponential in m .
K-Means Clustering
 $\min_{c_1, \dots, c_k} \sum_{i=1}^n \min_{j=1, \dots, k} \|x_i - c_j\|^2$
 1. Assign every data point to a cluster.
 2. Recalculate cluster means (Sum over all i in c_j / $|c_j|$)
 Repeat until not changing

Taylor Expansion
 $f(x) \approx f(a) + f'(a)(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f'''(a)}{3!}(x-a)^3 + \dots$
 $\frac{1}{1-x} = \sum_{n=0}^{\infty} x^n$, $\cos(x) = \sum_{n=0}^{\infty} \frac{(-1)^n x^{2n}}{(2n)!}$
 $\sin(x) = \sum_{n=0}^{\infty} \frac{(-1)^n x^{2n+1}}{(2n+1)!} = \sum_{n=1}^{\infty} \frac{(-1)^{n-1} x^{2n-1}}{(2n-1)!}$

Lagrange Multiplier for optimization
 Constraint $g(x) = 0$. $L(x, \lambda) = f(x) - \lambda g(x)$
 Lagrange multiplier must equal 0 \Rightarrow critical points.
 Check 2nd derivative or compare values to determine if maximal.
 Example: $\log \sum_{i=1}^n x_i$. Take care that $\sum_{i=1}^n x_i = 1$.
 $L(\theta, \lambda) = \sum_{i=1}^n x_i \log(x_i) + \lambda \left(\sum_{i=1}^n x_i - 1 \right)$
 $\frac{\partial L}{\partial x_i} = 0 \Rightarrow \theta_i = \dots$
 $\frac{\partial L}{\partial \lambda} = 0 \Rightarrow \theta_i = 1$
 From the constant $g=0$ we can now find λ .

Perceptron and SVM
 AA is always Kernel Requirements: Mat: $K(x, x')$ is also a Gram matrix.
 Pos. semi-def. $\rightarrow K(x, x') = k(x, x')$ Gram Matrix K must be pos. semi-definite. $k(x_0, x_0) \geq 0$.
 $k(x_0, x_0) = \sum_{i=1}^n x_{0i}^2 = \|x_0\|_2^2 \geq 0$
 Most also be symmetric and pos. semi-definite.
 So $k_{12} = k_{21}$ is not necessarily a kernel. $k_{12} = k_{21}$ is not necessarily a kernel.
 For det: $K \succeq 0$. EM Algorithm
 Have Data points but don't know their color or distribution. Assume Mixture Distribution. To guess some parameters. 2. Expect those PDFs of every data point to get likelihood of red/blue for data point given the guessed parameters.
 3. create weights $w = \frac{\text{red}}{\text{red} + \text{blue}}$ and send for blue.
 4. Compute mean and variance for red and blue PDF. \Rightarrow New params for red and blue PDF.
 estimate = $\text{sum}(\text{weights} * (\text{data} - \text{old Mean}))$
 estimateMean = $\text{sum}(\text{data} * \text{weights}) / \text{sum}(\text{weights})$
 Neural Network
 Given $\phi(z) = \frac{1}{1 + \exp(-z)}$ value, layer $\Rightarrow \phi(w_0 + \sum w_j x_j)$
 Treat final output as 1 if greater than 0.5, otherwise 0.
 Give weights for each output is not necessarily 0.5, otherwise 0.
 OR: $A \cdot B = \text{round}(\phi(w_0 + w_1 x_1 + w_2 x_2)) \Rightarrow w = [0.5, 1]^T$
 $A \cdot B = [0.5, 1]^T$ (for this ϕ)
 Task: Create a MLP with only one hidden layer and an arbitrary number of nodes that implements XORs. Possible approach: turn expr into ORs of ANDs and try patching all ORs into second layer.
 1st layer: $(A \cdot B) \Rightarrow h_1 = \phi(A + B \cdot C - 2.5)$ To get close enough to 1.
 $(A \cdot B) \Rightarrow h_2 = \phi(B \cdot C - D - 1.5)$ To get close enough to 0.
 $(A \cdot B) \Rightarrow h_3 = \phi(-A + B \cdot C - 0.4)$ To get close enough to 1.
 Task: XOR with as few units as possible.
 $\text{out} = \phi(-\phi(A+B) + 0.5 \cdot \phi(2A+2B))$
 $\text{out} = 0.84$

Taylor Expansion
 $f(x) \approx f(a) + f'(a)(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f'''(a)}{3!}(x-a)^3 + \dots$
 $\frac{1}{1-x} = \sum_{n=0}^{\infty} x^n$, $\cos(x) = \sum_{n=0}^{\infty} \frac{(-1)^n x^{2n}}{(2n)!}$
 $\sin(x) = \sum_{n=0}^{\infty} \frac{(-1)^n x^{2n+1}}{(2n+1)!} = \sum_{n=1}^{\infty} \frac{(-1)^{n-1} x^{2n-1}}{(2n-1)!}$

Lagrange Multiplier for optimization
 Constraint $g(x) = 0$. $L(x, \lambda) = f(x) - \lambda g(x)$
 Lagrange multiplier must equal 0 \Rightarrow critical points.
 Check 2nd derivative or compare values to determine if maximal.
 Example: $\log \sum_{i=1}^n x_i$. Take care that $\sum_{i=1}^n x_i = 1$.
 $L(\theta, \lambda) = \sum_{i=1}^n x_i \log(x_i) + \lambda \left(\sum_{i=1}^n x_i - 1 \right)$
 $\frac{\partial L}{\partial x_i} = 0 \Rightarrow \theta_i = \dots$
 $\frac{\partial L}{\partial \lambda} = 0 \Rightarrow \theta_i = 1$
 From the constant $g=0$ we can now find λ .

Generalization Error $E_{gen}(w) = \int_{\mathcal{X}} \ell(w; x, y) p(x, y) dx$

Expected Error of Prediction Function $f(x)$ is $I[f] = \int_{\mathcal{X}} \ell(f(x), y) p(x, y) dx$ where ℓ is a loss function and p is the unknown joint probability distribution for x and y . We don't know p so we cannot compute $I[f]$.

Instead, we compute the **empirical error** $I_n[f] = \frac{1}{n} \sum_{i=1}^n \ell(f(x_i), y_i)$. The generalization error $I[f]$ is said to be **empirical risk** $I_n[f]$.

An algorithm is said to be **empirical risk minimizing** if this goes towards 0 wrt. n .

Regularization (see also page 1)

Lasso: $\|w\|_1 = \sum w_i \rightarrow$ Euclidean drives quickly towards 0, softens ridge.

Ridge: $\|w\|_2 = \sqrt{\sum w_i^2}$

Euclidean: $\|w\|_2 = \sqrt{\sum w_i^2}$

Ridge vs Euclidean: \Rightarrow more stable near 0.

Proving positive semi-definiteness $X^T X \succeq 0$

- all Eigenvalues are ≥ 0
- $\forall v \in \mathbb{R}^n, v^T X^T X v = \|Xv\|_2^2 \geq 0$
- Gram matrix $X^T X$ was made from lin. indep. vectors x_i .
- Sylvester's Criterion: leading principal minors positive (\Rightarrow all \times $\neq 0$)
- left hand side principal minors $\neq 0$, upper triangular submatrix's determinant pos. definite iff all those are positive.
- \Rightarrow Test by reducing to upper triangular (row \rightarrow col) like in Gaussian Elimination, careful about pivots to keep sign of det correct! Non check that all diagonal elements are > 0 . (then ≥ 0)

Cumulative Distribution Function

COF = $F(x) = P(X \leq x)$ input

PDF = $f(x) = F'(x)$

Kernels Method

linear $(k(x, x)) = x^T x$

kernel $(k(x, x)) = (x^T x)^2$

$\phi(x) = \begin{pmatrix} x \\ x^2 \end{pmatrix}$, s.t. $k(x, x) = \phi(x)^T \phi(x) = x^2 + x^4$

$\int_{\mathcal{X}} f(x) p(x) dx = \int_{\mathcal{X}} f(x) p(x) dx$

Probability

$P(a, b|c) = P(a, b, c) / P(c)$

$P(a, b|c) = P(a, b, c) / P(c)$

$P(c|a, b) = P(a, b, c) / P(a, b)$

$P(c) = P(a, b, c) + P(a, \bar{b}, c) + P(\bar{a}, b, c) + P(\bar{a}, \bar{b}, c)$

EM-Atlap Concor

Goal: $P(X=x) = \sum_{i=1}^K \pi_i \delta_{x_i}$ if $x \in M$ possible

For each i , $\pi_i \geq 0$, $\sum \pi_i = 1$

EM-step: $\theta_{new} = \arg \max_{\theta} \sum_{i=1}^K \log \pi_i P(x_i; \theta)$ which can be written as $\sum_{i=1}^K \log \pi_i P(x_i; \theta)$

$\theta = \{ \mu, \sigma \}$

$\pi_i = \frac{1}{n} \sum_{j=1}^n \mathbb{1}_{\{x_j = x_i\}}$

EM-step: $\theta_{new} = \arg \max_{\theta} \sum_{i=1}^K \log \pi_i P(x_i; \theta)$

EM-Atlap Concor

Goal: $P(X=x) = \sum_{i=1}^K \pi_i \delta_{x_i}$ if $x \in M$ possible

For each i , $\pi_i \geq 0$, $\sum \pi_i = 1$

EM-step: $\theta_{new} = \arg \max_{\theta} \sum_{i=1}^K \log \pi_i P(x_i; \theta)$

$\theta = \{ \mu, \sigma \}$

$\pi_i = \frac{1}{n} \sum_{j=1}^n \mathbb{1}_{\{x_j = x_i\}}$

EM-step: $\theta_{new} = \arg \max_{\theta} \sum_{i=1}^K \log \pi_i P(x_i; \theta)$

EM-Atlap Concor

Goal: $P(X=x) = \sum_{i=1}^K \pi_i \delta_{x_i}$ if $x \in M$ possible

For each i , $\pi_i \geq 0$, $\sum \pi_i = 1$

EM-step: $\theta_{new} = \arg \max_{\theta} \sum_{i=1}^K \log \pi_i P(x_i; \theta)$

$\theta = \{ \mu, \sigma \}$

$\pi_i = \frac{1}{n} \sum_{j=1}^n \mathbb{1}_{\{x_j = x_i\}}$

EM-step: $\theta_{new} = \arg \max_{\theta} \sum_{i=1}^K \log \pi_i P(x_i; \theta)$

EM-Atlap Concor

Goal: $P(X=x) = \sum_{i=1}^K \pi_i \delta_{x_i}$ if $x \in M$ possible

For each i , $\pi_i \geq 0$, $\sum \pi_i = 1$

EM-step: $\theta_{new} = \arg \max_{\theta} \sum_{i=1}^K \log \pi_i P(x_i; \theta)$

$\theta = \{ \mu, \sigma \}$

$\pi_i = \frac{1}{n} \sum_{j=1}^n \mathbb{1}_{\{x_j = x_i\}}$

EM-step: $\theta_{new} = \arg \max_{\theta} \sum_{i=1}^K \log \pi_i P(x_i; \theta)$

EM-Atlap Concor

Goal: $P(X=x) = \sum_{i=1}^K \pi_i \delta_{x_i}$ if $x \in M$ possible

For each i , $\pi_i \geq 0$, $\sum \pi_i = 1$

EM-step: $\theta_{new} = \arg \max_{\theta} \sum_{i=1}^K \log \pi_i P(x_i; \theta)$

$\theta = \{ \mu, \sigma \}$

$\pi_i = \frac{1}{n} \sum_{j=1}^n \mathbb{1}_{\{x_j = x_i\}}$

EM-step: $\theta_{new} = \arg \max_{\theta} \sum_{i=1}^K \log \pi_i P(x_i; \theta)$

EM Slides (Soft) $\forall x_i \in \mathcal{X}$ second

E: Compute cluster membership probabilities $\gamma_i(x_j)$ for each i , given μ_j, σ_j from the previous step

$\gamma_i(x_j) = \frac{P(x_j | \mu_i, \sigma_i)}{\sum_{k=1}^K P(x_j | \mu_k, \sigma_k)}$

$\mu_j = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) x_i$

$\sigma_j^2 = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) (x_i - \mu_j)^2$

EM: $\mu_j = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) x_i$

$\sigma_j^2 = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) (x_i - \mu_j)^2$

EM Slides (Soft) $\forall x_i \in \mathcal{X}$ second

E: Compute cluster membership probabilities $\gamma_i(x_j)$ for each i , given μ_j, σ_j from the previous step

$\gamma_i(x_j) = \frac{P(x_j | \mu_i, \sigma_i)}{\sum_{k=1}^K P(x_j | \mu_k, \sigma_k)}$

$\mu_j = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) x_i$

$\sigma_j^2 = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) (x_i - \mu_j)^2$

EM: $\mu_j = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) x_i$

$\sigma_j^2 = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) (x_i - \mu_j)^2$

EM Slides (Soft) $\forall x_i \in \mathcal{X}$ second

E: Compute cluster membership probabilities $\gamma_i(x_j)$ for each i , given μ_j, σ_j from the previous step

$\gamma_i(x_j) = \frac{P(x_j | \mu_i, \sigma_i)}{\sum_{k=1}^K P(x_j | \mu_k, \sigma_k)}$

$\mu_j = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) x_i$

$\sigma_j^2 = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) (x_i - \mu_j)^2$

EM: $\mu_j = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) x_i$

$\sigma_j^2 = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) (x_i - \mu_j)^2$

EM Slides (Soft) $\forall x_i \in \mathcal{X}$ second

E: Compute cluster membership probabilities $\gamma_i(x_j)$ for each i , given μ_j, σ_j from the previous step

$\gamma_i(x_j) = \frac{P(x_j | \mu_i, \sigma_i)}{\sum_{k=1}^K P(x_j | \mu_k, \sigma_k)}$

$\mu_j = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) x_i$

$\sigma_j^2 = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) (x_i - \mu_j)^2$

EM: $\mu_j = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) x_i$

$\sigma_j^2 = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) (x_i - \mu_j)^2$

EM Slides (Soft) $\forall x_i \in \mathcal{X}$ second

E: Compute cluster membership probabilities $\gamma_i(x_j)$ for each i , given μ_j, σ_j from the previous step

$\gamma_i(x_j) = \frac{P(x_j | \mu_i, \sigma_i)}{\sum_{k=1}^K P(x_j | \mu_k, \sigma_k)}$

$\mu_j = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) x_i$

$\sigma_j^2 = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) (x_i - \mu_j)^2$

EM: $\mu_j = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) x_i$

$\sigma_j^2 = \frac{1}{n} \sum_{i=1}^n \gamma_i(x_j) (x_i - \mu_j)^2$

More Probability

$P(w|x, y) = \frac{P(w, x, y)}{P(x, y)}$

$P(w, x, y) = P(w|x, y) P(x, y)$

$P(w|x, y) = \frac{P(w, x, y)}{P(x, y)}$

$P(w|x, y) = \frac{P(w, x, y)}{P(x, y)}$

More Probability

$P(w|x, y) = \frac{P(w, x, y)}{P(x, y)}$

$P(w, x, y) = P(w|x, y) P(x, y)$

$P(w|x, y) = \frac{P(w, x, y)}{P(x, y)}$

$P(w|x, y) = \frac{P(w, x, y)}{P(x, y)}$

More Probability

$P(w|x, y) = \frac{P(w, x, y)}{P(x, y)}$

$P(w, x, y) = P(w|x, y) P(x, y)$

$P(w|x, y) = \frac{P(w, x, y)}{P(x, y)}$

$P(w|x, y) = \frac{P(w, x, y)}{P(x, y)}$

More Probability

$P(w|x, y) = \frac{P(w, x, y)}{P(x, y)}$

$P(w, x, y) = P(w|x, y) P(x, y)$

$P(w|x, y) = \frac{P(w, x, y)}{P(x, y)}$

$P(w|x, y) = \frac{P(w, x, y)}{P(x, y)}$

More Probability

$P(w|x, y) = \frac{P(w, x, y)}{P(x, y)}$

$P(w, x, y) = P(w|x, y) P(x, y)$

$P(w|x, y) = \frac{P(w, x, y)}{P(x, y)}$

$P(w|x, y) = \frac{P(w, x, y)}{P(x, y)}$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$

Normal Distribution

$f(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Expected Value

$E[x] = \int_{-\infty}^{\infty} x f(x) dx = \mu$

$E[x^2] = \int_{-\infty}^{\infty} x^2 f(x) dx = \mu^2 + \sigma^2$

$E[x^3] = \int_{-\infty}^{\infty} x^3 f(x) dx = \mu^3 + 3\mu\sigma^2$

$E[x^4] = \int_{-\infty}^{\infty} x^4 f(x) dx = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$