

$$\|\mathbf{x}^{(a)} - \mathbf{x}^{(b)}\|_2 = \sqrt{\sum_{j=1}^D (x_j^{(a)} - x_j^{(b)})^2} \quad \text{cosine}(\mathbf{x}^{(a)}, \mathbf{x}^{(b)}) = \frac{\mathbf{x}^{(a)} \cdot \mathbf{x}^{(b)}}{\|\mathbf{x}^{(a)}\|_2 \|\mathbf{x}^{(b)}\|_2}$$

$$\text{Accuracy} = \frac{\text{Num Correct Predictions}}{\text{Num Predictions}} = \frac{\sum_{i=1}^N \mathbb{I}[t^{(i)} = y^{(i)}]}{N}$$

$$H(X) = -\mathbb{E}_{X \sim p}[\log_2 p(X)] = -\sum_{x \in X} p(x) \log_2 p(x)$$

$$H(X, Y) = - \sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(x, y)$$

$$H(Y|X = x) = - \sum p(y|x) \log_2 p(y|x) \quad p(y|x) = p(x,y)$$

$$p(x) = \sum_y p(x, y)$$

$H(Y|Y) = 0$. $H(Y|X) \leq H(Y)$ **Exp. Cond. Entro:** $H(Y|X) = \mathbb{E}_x[H(Y|X)] = -\sum_{x \in X} \sum_{y \in Y} p(x, y) \log_2 p(y|x)$ **Information gain:** $IG(Y|X) = H(Y) - H(Y|X)$

(bias goes in w)
$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$
 Squared loss function (bad x):

$\mathcal{L}(y, t) = \frac{1}{2}(y - t)^2$ $\mathcal{E}(w, b) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y^{(i)}, t^{(i)})$

$$\mathcal{E}(\mathbf{w}) = \frac{1}{2N} \|\mathbf{X}\mathbf{w} - \mathbf{t}\|_2^2 \quad \nabla_{\mathbf{w}} \mathcal{E}(\mathbf{w}) = \frac{1}{N} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{t})$$

$\underline{w} = (X^T X)^{-1} X^T t$ Gradient Desc: $\underline{w} \leftarrow \underline{w} - \alpha \nabla_{\underline{w}} \mathcal{E}$ Polynomial

$$\mathcal{R}(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_2^2 = \frac{1}{2} \sum_j w_j^2$$

$$z = \mathbf{w}^T \mathbf{x}$$

$$\therefore \sigma(\underset{\uparrow}{z}) = \frac{1}{1 + e^{-z}}$$

Cross Entropy loss:

$$= -t \log(y) - (1 - t) \log(1 - y)$$

$$\frac{\partial \mathcal{L}}{\partial w_j} = \frac{\partial \mathcal{L}}{\partial y} \cdot y(1-y) \cdot x_j = \left(-\frac{t}{y} + \frac{1-t}{1-y}\right) \cdot y(1-y) \cdot x_j = (y-t)x_j$$

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{\alpha}{N} \sum_{i=1}^N (y^{(i)} - t^{(i)}) \mathbf{x}^{(i)}$$

Fancy words: (non) Parametric, Hypothesis (space), inductive bias.

$$\frac{\partial \mathcal{E}}{\partial \mathbf{w}} = \frac{1}{N} \sum_{i=1}^N \frac{\partial \mathcal{L}(\mathbf{y}^{(i)}, \mathbf{t}^{(i)})}{\partial \mathbf{w}}$$

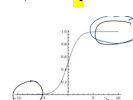
```
mini batches of size k##for each mini batch #estimate the gradient using the mini batch size
# update the parameters based on the estimate. Each pass of the inner loop is an iteration
-> One iteration is one update for each weight Each pass in the outerloop is called an epoch:
-> one epoch is one pass over the entire data-set
num_iter = num_epoch * N/k
```

$$z = Wx + b \quad y = \text{softmax}(z) \quad \mathcal{L}_{CE} = -\mathbf{t}^T \log(\mathbf{y})$$

$$y_k = \text{softmax}(z_1, \dots, z_K)_k = \frac{e^{z_k}}{\sum_{m=1}^K e^{z_m}} \quad \frac{\partial \mathcal{L}_{CE}}{\partial w_k} = \left(\frac{\partial \mathcal{L}_{CE}}{\partial z_k} \right) \left(\frac{\partial z_k}{\partial w_k} \right) = (y_k - t_k) \cdot x$$

$$\mathbf{w}_k \leftarrow \mathbf{w}_k - \frac{\alpha}{N} \sum_{i=1}^N (y_k^{(i)} - t_k^{(i)}) \cdot \mathbf{x}^{(i)}$$

$$y = \phi(\mathbf{w}^\top \mathbf{x} + b)$$



$$\text{ReLU}(z) = \max(z, 0)$$

$$\begin{aligned} \mathbf{h}^{(1)} &= \mathbf{f}(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)}) \\ \mathbf{h}^{(2)} &= \mathbf{f}(\mathbf{W}^{(2)} \mathbf{h}^{(1)} + \mathbf{b}^{(2)}) \\ &\vdots \\ \mathbf{h}^{(L-1)} &= \mathbf{f}(\mathbf{W}^{(L-1)} \mathbf{h}^{(L-2)} + \mathbf{b}^{(L-1)}) \\ \mathbf{z} &= \mathbf{W}^{(L)} \mathbf{h}^{(L-1)} + \mathbf{b}^{(L)} \\ \mathbf{y} &= \text{softmax}(\mathbf{z}) \end{aligned}$$

```

for i = 1:N
    Compute v_i as function of P_A(v_i)
     $\overline{v_N} = i$ 
    for i = N-1:1
         $\overline{v_i} = \sum_{j \in CH(v_i)} v_j \frac{\partial v_j}{\partial v_i}$ 
    end
end

```

For $z = W_1x + b$, $h = \sigma(z)$, $y = W_2h + b_2$, $L = \frac{1}{2} \|t - y\|^2$,
 $L = \frac{1}{2} \cdot \overline{W_1}(y - t)$, $\overline{W_2} = \overline{y}h^{-1}$, $b_2 = \overline{y}$, $h = W_2^{-1}\overline{y}$, $\overline{z} = \overline{h} \circ \sigma'(z)$
 $\overline{W_1} = \frac{\overline{z}}{x}$, $\overline{b_1} = \overline{z}$. In sigmoid activation with squared loss,
Want $\frac{\partial L}{\partial w} \cdot \frac{\partial L}{\partial b}$.

$$L = \frac{1}{2}(y - t)^2, y = \sigma(z), z = wx + b$$

$$\frac{\partial L}{\partial y} = y - t, \quad \frac{\partial L}{\partial z} = \frac{\partial L}{\partial y} \sigma'(z), \quad \frac{\partial L}{\partial w} = \frac{\partial L}{\partial z} x, \quad \frac{\partial L}{\partial b} = \frac{\partial L}{\partial z}$$

$\bar{\mathcal{L}} = 1$	$\bar{\mathcal{L}} = 1$
$\bar{y}_k = \mathcal{L}(y_k - t_k)$	$\bar{\mathbf{y}} = \bar{\mathcal{L}}(\mathbf{y} - \mathbf{t})$
$w_{ki}^{(2)} = \bar{y}_k h_i$	$\bar{\mathbf{W}}^{(2)} = \bar{\mathbf{y}} \mathbf{h}^T$
$b_k^{(2)} = \bar{y}_k$	$\bar{\mathbf{b}}^{(2)} = \bar{\mathbf{y}}$
$h_i = \sum_k \bar{y}_k w_{ki}^{(2)}$	$\bar{\mathbf{h}} = \bar{\mathbf{W}}^{(2)T} \bar{\mathbf{y}}$
$\mathbf{z}_i = h_i \sigma(\mathbf{z}_i)$	$\bar{\mathbf{z}} = \mathbf{h} \circ \sigma^T(\mathbf{z})$
$w_{ij}^{(1)} = \mathbf{z}_i x_j$	$\bar{\mathbf{W}}^{(1)} = \bar{\mathbf{z}} \mathbf{x}^T$
$b_j^{(1)} = \mathbf{z}_j$	$\bar{\mathbf{b}}^{(1)} = \bar{\mathbf{z}}$

$$\mathbf{z} \in \mathcal{R}^N, \mathbf{y} \in \mathcal{R}^M \quad \bar{z}_j = \sum_k \bar{y}_k \frac{\partial y_k}{\partial z_j} \quad \bar{\mathbf{z}} = \frac{\partial \mathbf{y}^\top}{\partial \mathbf{z}} \bar{\mathbf{y}},$$

$$\left(\frac{\partial \mathbf{y}}{\partial \mathbf{z}}\right)_{M \times N} = \begin{pmatrix} \frac{\partial y_1}{\partial z_1} & \dots & \frac{\partial y_1}{\partial z_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial z_1} & \dots & \frac{\partial y_m}{\partial z_n} \end{pmatrix}$$

Bias Variance decomp:

Variance: $\text{Var}(y)$, How the average prediction (across various training set samples) differs from the prediction we get from one particular training set. High capacity models have high variance.

Bayes Error: $\text{var}(t)$ The inherent unpredictability of the targets.

$$Var(1/m \sum_{i=1}^m h_i(x)) = (1/m(1-p) + \rho)\sigma^2$$

May be more advantageous to introduce more variability into algo, as long as it reduces the correlation between samples. Bagging reduces overfitting by averaging predictions but does not reduce bias. **Random Forests:** Bagged decision trees where when choosing each of node, choose split on random subset from d input features. Improves the variance reduction of bagging by reduced correlation b/w tree - p . Bagging does not control Bayes error.

$$\sigma = \sqrt{\frac{\sum (X - \mu)^2}{N}}$$

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

$$\mathbb{E}[(\hat{f}(\mathbf{x}) - t)^2] = \underbrace{(\mathbb{E}[y_* - y])^2}_{\text{Bias}^2} + \underbrace{\text{Var}(y)}_{\text{Noise}} + \underbrace{\text{Var}(t)}_{\text{Model Variance}}$$

Naive Bayes

$$p(\mathbf{x}) = \sum_c p(\mathbf{x}|c)p(c)$$

binary variables requires $2^{D+1} - 1$ **Solution:** assume each word feature x_i is conditionally independent given class c , $2D + 1$ now

Naive Bayes: Model $P(x, y)$. Know $P(x|y) \rightarrow P(x|y) = \frac{P(x, y)}{P(y)}$
 Know $P(x|y) \rightarrow P(x|y) = \frac{P(x, y)}{P(y)}$. Let Bern, $\theta \in [0, 1]$.
 $P(x, y) = \theta \Rightarrow P(x, y) = \theta \Rightarrow P(x|y) = \theta^{x(1-\theta)^{1-x}}$ (weighted)
 Assume $\{x_1, \dots, x_N\}$ i.i.d. Bern $\Rightarrow P(x_1, \dots, x_N) = \prod_{i=1}^N \theta^{x_i(1-\theta)^{1-x_i}} = \theta^{\sum_{i=1}^N x_i} (1-\theta)^{N - \sum_{i=1}^N x_i}$
 $\log(L(\theta)) = \sum_{i=1}^N \{x_i \log \theta + (1-x_i) \log(1-\theta)\}$. $\hat{\theta}_{MLE} = \arg \max_{\theta} L(\theta)$.
 (or) Lib example: $\hat{\theta}_{MLE} = \frac{N_H + N_F}{N_H + N_F + 1}$. $\hat{\theta}_{MLE} = \arg \max_{\theta} L(\theta)$
 X is Bow features: $P(x|c) = \frac{P(x, c)}{P(c)}$
 Here $P(x) = \sum_{c \in \text{category}} P(x, c) P(c)$ (class weighted: $P(c=1) = \pi$)
 $P(c=j|1, c) = \theta_j c^j$. 2D + 1 param not 2^{d+1}
 Naive Bayes log likelihood: $\theta = \{ \theta_1, \theta_2 \}$ $L(\theta) = \sum_{i=1}^N \log P(c^{(i)} | x^{(i)})$
 $= \sum_{i=1}^N \{ \log \theta_1^{x_i} \theta_2^{1-x_i} \} = \sum_{i=1}^N \{ x_i \log \theta_1 + (1-x_i) \log \theta_2 \}$
 Optimizer: $\pi = \frac{\sum_{i=1}^N \mathbb{I}(c^{(i)} = 1)}{N}$ (Bern LL: $\theta_1 = \pi$)
 Optimizer: $\theta_c = \frac{\sum_{i=1}^N \mathbb{I}(x^{(i)} = 1 \wedge c^{(i)} = c)}{\sum_{i=1}^N \mathbb{I}(c^{(i)} = c)}$ (Bern LL: $\theta_2 = \frac{\sum_{i=1}^N x_i \mathbb{I}(c^{(i)} = c)}{\sum_{i=1}^N \mathbb{I}(c^{(i)} = c)}$)
 Total data \Rightarrow overall data Bern LL: $\theta = \{ \theta_1, \theta_2 \}$. So $P(\theta|D) \propto P(\theta) P(D|\theta)$
 Beta dist: $P(\theta; a, b) \propto \theta^a (1-\theta)^b$

Given data: $p(\theta, D)$ and $p(D, \theta)$
 $\hat{\theta}_{MAP} = \arg \max_{\theta} p(\theta | D) = \arg \max_{\theta} p(\theta, D) = \arg \max_{\theta} p(\theta) p(D | \theta)$
 $= \arg \max_{\theta} (\log p(\theta) + \log p(D | \theta))$
 $\hat{\theta}_{MAP} = \frac{N_H + a - 1}{N_H + N_L + a + b - 2}$

Bayes's Classifier: Given features $x = [x_1, x_2, \dots, x_D]^T$ we want to compute class probabilities using Bayes Rule $P(c|x) = \frac{p(c)p(x|c)}{p(x)}$. At training time, estimate parameters using log likelihood, compute co-occurrence counts of each feature with the labels. At test time, apply Bayes' Rule. **MLE** Want to maximize the log likelihood: $l(\theta) = \sum \log(p(x_i|\theta))$. Too little data can cause overfitting.

Bayesian Parameter Estimation: MLE treat observation as random variables, but parameters are not. Bayesian treats parameters as random variables as well. Define *prior distribution* $p(\theta)$, which encodes beliefs about parameter before observing data. *likelihood* $p(D|\theta)$, same as MLE. Updating beliefs by computing posterior distribution using Bayes' Rule:

$$p(\theta|D) = \frac{p(\theta)p(D|\theta)}{\int p(\theta')p(D|\theta')d\theta'}$$

MAP Estimator: Find the most likely parameter settings under the posterior. Want to maximize the log likelihood: $l(\theta) = \log(p(\theta)) + \sum \log p(x_i|\theta)$.

Exercise: Learning each θ_{jk}

Formula

$$\hat{\theta}_{ML} = \frac{N_{Hj}}{N_{Hj} + N_T}$$

$$E[\theta|D] = \frac{N_{Hj} + a}{N_{Hj} + N_T + a + b}$$

$$\hat{\theta}_{MAP} = \frac{N_{Hj} + a - 1}{N_{Hj} + N_T + a + b - 2}$$

Exercise Compute the MAP estimate for π and each θ_{jk} in this data set, assuming a Beta distribution priors for each parameter with $a=2, b=2$.

- he has funny and sad, "positive"
- movie was sad, "positive"
- great movie, "positive"
- movie was not funny, "negative"
- movie was not sad, "negative"

Then, use your parameters to estimate the probabilities that these reviews are positive:

- movie was funny
- he was not funny

Gaussian Discriminant Analysis:

PDF of the univariate Gaussian distribution:

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

PDF of the multivariate Gaussian distribution:

$$\mathcal{N}(x; \mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right]$$

$$\text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N} \quad \Sigma = \text{Cov}(x) = E[(x - \mu)(x - \mu)^T]$$

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x^{(i)} \quad \hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \hat{\mu})(x^{(i)} - \hat{\mu})^T$$

Gaussian Discriminative Analysis (GDA):

Assumes that $p(x|t)$ is distributed according to multivariate Gaussian distribution.

$$p(x|t = k) = \frac{1}{(2\pi)^D |\Sigma_k|^2} \exp\left[-\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)\right]$$

Where $|\Sigma_k|$ denotes the determinant of the matrix, and D is dimension of x . Each class k has a mean vector μ_k and a covariance matrix Σ_k . Use MLE to find empirical mean and empirical covariance matrix: $\hat{\mu}_k$ and $\hat{\Sigma}_k$. Σ_k has $O(D^2)$ parameters - could be hard to estimate

To find log-likelihood of GDA: Assume data x are drawn from gaussian, with μ and Σ . We calculate log-likelihood $l(\mu, \Sigma)$, then take derivative w.r.t μ and Σ , evaluate at zero to find the MLE $\hat{\mu}$ and $\hat{\Sigma}$.

Decision boundary of GDA: usually quadratic. It is based on class posterior. GDA makes decisions by comparing class probabilities.

$$\log p(t_k|x) = \log p(x|t_k) + \log p(t_k) \quad \text{Bayes rule}$$

$$= -\frac{d}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma_k^{-1}| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log p(t_k) - \log p(x)$$

Now, suppose I have two classes, l and k . Then the decision boundary equals the follow: decision boundary is $\log p(t_k|x) = \log p(t_l|x)$, which is equivalent to

$$(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) = (x - \mu_l)^T \Sigma_l^{-1} (x - \mu_l) + C_{kl}$$

Quadratic function in $x \Rightarrow$ **quadratic decision boundary!**

What is C_{kl} ? What if $\Sigma_k = \Sigma_l$? \Rightarrow can be derived. Easy.

Spectral Decomposition

For matrix A , anything satisfying $Av = \lambda v$, is an eigenvector and its associated eigenvalue. Our $D \times D$ matrix will have at most D eigenvalues.

Since our A is symmetric then the Spectral Theorem tells us that: All its **eigenvalues** are **real-valued**.

There is a full set of linearly independent **eigenvectors** that form a basis for \mathbb{R}^D .

These **eigenvectors** can be chosen to be **real-valued** and **orthonormal**.

Quadratic Forms of Symmetric matrices

$v^T A v > 0$, for all $v \neq 0 \rightarrow$ the quadratic form curves upwards; A is **positive definite**; $A \succ 0$

$v^T A v \geq 0$ for all $v \rightarrow A$ is **positive semidefinite**; $A \succcurlyeq 0$

$v^T A v < 0$ for all $v \neq 0$, we say A is **negative definite**, denoted $A \prec 0$

$v^T A v$ can be positive or negative then A is **indefinite**

Claim: A matrix A is positive semi-definite iff all of its eigenvalues are nonnegative.

If A is both diagonal and positive definite (i.e. its diagonal entries are positive), then the ellipses are axis-aligned.

Comparing GDA vs Logistic Regression: GDA makes stronger modeling assumption: assumes class-conditional data is multivariate Gaussian. If this is true, GDA is asymptotically efficient.

But LR is more robust, less sensitive to incorrect modeling assumptions. When the class-conditional distributions are non-Gaussian (true almost always), LR usually beats GDA.

Also, GDA has quadratic (when not shared covariance), LR has linear decision boundary.

Gaussian Naive Bayes classifier assumes that the likelihoods are Gaussian:

$$p(x_j | c = k) = \frac{1}{\sqrt{2\pi}\sigma_{jk}} \exp\left[-\frac{(x_j - \mu_{jk})^2}{2\sigma_{jk}^2}\right]$$

- this is just a 1-dim Gaussian, one for each input dimension
- Model the same as GDA with diagonal covariance matrix

Unsupervised learning

K-means: Group data points into clusters. Relies on assumption that data depends on latent variables. The objective:

$$\min_{\{m_k\}, \{r^{(n)}\}} \sum_{n=1}^N \sum_{k=1}^K r_k^{(n)} \|x^{(n)} - m_k\|^2$$

Where $\{m_k\}$ are the k means, $r^{(n)} \in \mathbb{R}^K$ is the assignment vector for $x^{(n)}$. $r^{(n)}$ is one-hot if hard k-means and like a probability distribution if soft k-means. Finding the min is NP-hard in general.

Alternating minimization involves fixing one of $\{m_k\}$, or $\{r^{(n)}\}$, optimizing on that, then alternate to other. First initialize clusters, then iterate between assignment step (assign data points to nearest clusters), and refitting step (move every cluster centre to the means of the data points assigned to it). This algo is guaranteed to converge to a local optimum. **Test of convergence:** if the assignments don't change over an iteration, then stop. **Refitting Step:** For each centre:

$$m_k = (\sum_n r_k^{(n)} x^{(n)}) / (\sum_n r_k^{(n)})$$

Assignment Step: For soft k-means:

$$r_k^{(n)} = \frac{\exp(-\beta \|m_k - x^{(n)}\|^2)}{(\sum_j \exp(-\beta \|m_j - x^{(n)}\|^2))}$$

$$r^{(n)} = \text{softmax}(-\beta \|m_k - x^{(n)}\|^2)_{k=1}^K$$

For regular: $r_k^{(n)} = \mathbb{I}[k^{(n)} = k]$, for $k = 1, \dots, K$, where $k^{(n)} = \text{argmin}_k \|m_k - x^{(n)}\|^2$. As β goes to infinity, soft k-means becomes regular.

By marginalizing over z , we get a density over the observables:

$$p(x) = \sum_z p(x, z) = \sum_z p(x|z) p(z)$$

This is called a **latent variable model**.

If $p(z)$ is a categorical distribution, this is a **mixture model**, and different values of z correspond to different **components**.

GMM:

$$p(x) = \sum_{k=1}^K \pi_k \mathcal{N}(x | \mu_k, \Sigma_k) \quad \sum_{k=1}^K \pi_k = 1$$

Maximum likelihood principle:

$$\log p(D; \theta) = \sum_{n=1}^N \log p(x^{(n)}; \theta)$$

Gaussian Mixture Models: A generative approach to clustering. We need a method to judge what it means to cluster the data well. Assume that the data was generated from a model. We then adjust the model parameters to maximize the probability that the data was generated by our model. Given D , we assume that the data point, x , was generated by choosing a cluster $k \in \{1, \dots, K\}$, such that $p(z = k) = \pi_k$, and sampling from $\mathcal{N}(x | \mu_k, \Sigma_k)$, i.e. $p(x|z = k) = \mathcal{N}(x | \mu_k, \Sigma_k)$. To find the best parameters, we want to max:

$$\log p(D) = \sum_{n=1}^N \log p(x^{(n)}) = \sum_{n=1}^N \log \left(\sum_{k=1}^K \pi_k \mathcal{N}(x^{(n)} | \mu_k, \Sigma_k) \right)$$

EM Algorithm for GMM: A local algorithm for GMM. First initialize the starting μ_k and Σ_k . Then iterate between the E step and M step until convergence (check by looking at log likelihood).

E step: posterior probability over z over current model (how much we think a datapoint was generated by a cluster):

$$r_k^{(n)} = p(z^{(n)} = k | x^{(n)}) = \frac{\pi_k \mathcal{N}(x^{(n)} | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(x^{(n)} | \mu_j, \Sigma_j)}$$

M step: update the model parameters μ_k, Σ_k to optimize the expected complete data log-likelihood. $\mu_k = (1/N_k) \sum_{n=1}^N r_k^{(n)} x^{(n)}$, $\Sigma_k = N_k/N$, where $N_k = \sum_{n=1}^N r_k^{(n)}$.

$$\Sigma_k = \frac{1}{\sum_{i=1}^N r_k^{(i)}} \sum_{i=1}^N r_k^{(i)} (x^{(i)} - \mu_k)(x^{(i)} - \mu_k)^T$$

Both k-means and EM suffer

from getting stuck at local minima **BMM:** Bernoulli Mixture Model

$$p(x) = \prod_k \prod_j p(x_{ij} | z = k) p(z = k) \quad r_k^{(i)} = p(z^{(i)} = k | x^{(i)}) = \frac{p(z = k) p(x | z = k)}{\sum_{\ell} p(z = \ell) p(x | z = \ell)}$$

$$\text{Proj}_S(x) = \sum_{i=1}^K z_i u_i \quad \text{where } z_i = x^T u_i \quad \text{Proj}_S(x) = U z \quad \text{where } z = U^T (x - \hat{\mu})$$

Min reconstruction error/max variance of z 's

$$\min_U \frac{1}{N} \sum_{i=1}^N \|x^{(i)} - \hat{x}^{(i)}\|^2 \quad \max_U \frac{1}{N} \sum_{i=1}^N \|\hat{x}^{(i)} - \hat{\mu}\|^2$$

Because the columns of U are orthogonal, $U^T U = I$, the norm of centered reconstruction is equal to norm of representation:

$$\|\hat{x} - \hat{\mu}\|^2 = \|U z\|^2 = z^T U^T U z = z^T z = \|z\|^2$$

So, we need to choose $a = Q^T u$ that maximizes

Rotate X using $Q \rightarrow$ covariance of transformed data will be diagonal

Projected variance = constant - reconstruction error

PCA: ui perp basis

$$\hat{x} \approx \text{Proj}_S(x) = \sum_{i=1}^K z_i u_i \quad \text{where } z_i = x^T u_i \in \mathbb{R}$$

$$a^T \Lambda a = \sum_{j=1}^D \lambda_j a_j^2$$

Our goal is to show that $\frac{1}{N} \sum_{i=1}^N \|x^{(i)} - \hat{x}^{(i)}\|^2 = \text{const} - \frac{1}{N} \sum_{i=1}^N \|\hat{x} - \hat{\mu}\|^2$

Claim: The vectors $\hat{x}^{(i)} - \hat{\mu}$ and $\hat{x}^{(i)} - x^{(i)}$ are **orthogonal** (and we can thus apply the Pythagorean theorem).

Proof: First, we can write $\hat{x}^{(i)} = \hat{\mu} + U U^T (x^{(i)} - \hat{\mu})$.

We know that two vectors a, b are orthogonal $\Leftrightarrow a^T b = 0$.

So we compute $(\hat{x}^{(i)} - \hat{\mu})^T (\hat{x}^{(i)} - x^{(i)})$:

$$\begin{aligned} & (\hat{x}^{(i)} - \hat{\mu})^T (\hat{x}^{(i)} - x^{(i)}) \\ &= (x^{(i)} - \hat{\mu})^T U U^T (\hat{\mu} - x^{(i)} + U U^T (x^{(i)} - \hat{\mu})) \\ &= (x^{(i)} - \hat{\mu})^T U U^T (\hat{\mu} - x^{(i)}) + (x^{(i)} - \hat{\mu})^T U U^T (x^{(i)} - \hat{\mu}) \\ &= 0 \end{aligned}$$

The Pythagorean Theorem tells us:

$$\|\hat{x}^{(i)} - \hat{\mu}\|^2 + \|x^{(i)} - \hat{x}^{(i)}\|^2 = \|x^{(i)} - \hat{\mu}\|^2 \quad \text{for each } i$$

By averaging over data and from observance, we obtain

$$\begin{aligned} & \frac{1}{N} \sum_{i=1}^N \|\hat{x}^{(i)} - \hat{\mu}\|^2 + \frac{1}{N} \sum_{i=1}^N \|x^{(i)} - \hat{x}^{(i)}\|^2 \\ &= \underbrace{\frac{1}{N} \sum_{i=1}^N \|\hat{x}^{(i)} - \hat{\mu}\|^2}_{\text{projected variance}} + \underbrace{\frac{1}{N} \sum_{i=1}^N \|x^{(i)} - \hat{x}^{(i)}\|^2}_{\text{reconstruction error}} \\ &= \underbrace{\frac{1}{N} \sum_{i=1}^N \|x^{(i)} - \hat{\mu}\|^2}_{\text{constant}} \end{aligned}$$

Therefore, **projected variance = constant - reconstruction error**

Maximizing the variance is equivalent to minimizing the reconstruction error!

Principle Component Analysis (PCA): Reduces the dimension of data which can help with saving computation/memory and reducing overfitting.

Given datapoints, consider the sample mean $\hat{\mu}$ and the sample covariance $\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \hat{\mu})(x^{(i)} - \hat{\mu})^T$. Can find a subspace S of \mathbb{R}^D with orthonormal basis u_1, \dots, u_K that "best represents" our $x^{(i)} - \hat{\mu}$ (or minimizes reconstruction squared error loss or maximizes variance of reconstructions (sum of squares of distances to the mean)) by doing spectral decomposition of $\hat{\Sigma}$, sorting the eigenvalues (which are all non-negative as $\hat{\Sigma}$ is positive semidefinite) from biggest to smallest, then taking the first K eigenvectors as u_1, \dots, u_K . We project points $x^{(i)} - \hat{\mu}$ onto S by applying the matrix U^T for $U = (u_1, \dots, u_K)$. We can then

"unproject" by applying U . The features of the projected datapoints are uncorrelated.

Deriving PCA: Maximize $a^T \Lambda a = \sum_{j=1}^D \lambda_j a_j^2$ for $a = Q^T u$, assume u are in descending order, by observation, since u is a unit vector, then by unitarity, a is also a unit vector: $a^T a = u^T Q Q^T u = u^T u$. By inspection, set $a = \pm 1$, and $a_j = 0$ for $j \neq 1$. Hence, $u = Q a = q_1$ (the top eigenvector)