

1. SVM, $y_i(w \cdot x_i + b) \geq 1 - \xi_i, \xi_i \geq 0$. Slack variables, distance that need to move bad point. Learning rate $\sum \eta(i) = \infty$, but $\sum \eta(i)^2 < \infty$. e.g. $\eta(i) = \frac{c}{i}$. Momentum: $A_i = \vec{\theta}^{i+1} - \vec{\theta}^i = \eta \nabla f + \alpha A_{i-1}$.

2. Convex Set D . $x \in D, y \in D, \lambda x + (1-\lambda)y \in D, \lambda \in [0,1]$. Convex Function, with domain D . $f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2)$. $f(y) = f(x) + \nabla f(x)(y-x)$. Chain Rule = $f(x) = g(h(x)), \frac{\partial f}{\partial x} = \frac{\partial h(x)}{\partial x} \cdot \frac{\partial g}{\partial h}$. $f(x) = g(h(x)) = g\left(\begin{matrix} h_1(x) \\ h_2(x) \\ \vdots \\ h_k(x) \end{matrix}\right) \frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial h_1}{\partial x_1} & \frac{\partial h_1}{\partial x_2} & \dots & \frac{\partial h_1}{\partial x_n} \\ \frac{\partial h_2}{\partial x_1} & \frac{\partial h_2}{\partial x_2} & \dots & \frac{\partial h_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial h_k}{\partial x_1} & \frac{\partial h_k}{\partial x_2} & \dots & \frac{\partial h_k}{\partial x_n} \end{bmatrix} \begin{bmatrix} \frac{\partial g}{\partial h_1} \\ \frac{\partial g}{\partial h_2} \\ \vdots \\ \frac{\partial g}{\partial h_k} \end{bmatrix}$. Dual = $\max_{\alpha} \min_x D(x)$ for inequality, $\vec{\alpha} \geq 0$. KKT condition: $\min f(x)$, subject to $g_i(x) \leq 0, h_j(x) = 0$. $\frac{\partial L(x^*, \alpha^*, \beta^*)}{\partial x} = 0$ ① $h(x^*) = 0$ ② $\alpha_i^* g_i(x^*) = 0$ ③ Note in ①, in augmented form, do not take augmented.

3. $f(\vec{x})$ vector valued, $= [x_1^2, x_2^2, \dots, x_n^2]^T$. $\frac{\partial f}{\partial x} = \nabla f(\vec{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$. $\nabla_x f = \frac{\partial f}{\partial x_j} \leftarrow \text{column}$ $\frac{\partial f}{\partial x_j} \leftarrow \text{row}$. $f(\vec{x}) = \begin{bmatrix} f_1(\vec{x}) \\ f_2(\vec{x}) \end{bmatrix}, \nabla_x f = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} \end{bmatrix}$. ④ $g_i(x^*) \leq 0$ ⑤ $\alpha_i^* \geq 0$. Slater condition, exist $x, h(x) = 0, g(x) < 0$.

4. $\nabla_x^2 f$ is p.s.d. $\Leftrightarrow f$ is convex. positive definite (p.d.) \Leftrightarrow strictly convex. 12. Check if Dual applied. ① $f(x)$ is convex ② (inequality) $g(x)$ is convex ③ $h(x)$ is linear ④ Slater condition, exist $x, h(x) = 0, g(x) < 0$.

5. A is p.s.d if $x^T A x \geq 0$ for all x . $x^T A x = x_1^2 a_{11} + x_1 x_2 a_{12} + x_2 x_1 a_{21} + x_2^2 a_{22}$. A is p.d if $x^T A x > 0$ for all x . 13. ADMM - P1 $\arg \min_{x, z} \frac{1}{2} (z - u)^2 + \lambda |x| + \alpha (x - z)$ original problem $\arg \min \frac{1}{2} \|x - u\|_2^2 + \lambda \|x\|_1$. ① $\min \lambda |x| + \alpha x, (\lambda > 0)$ regularization constant. if $\alpha > \lambda$, then $-\infty$, because set $x = -\infty$. $\alpha = \lambda$, then 0, because $x = 0$ or $x < 0 \Rightarrow \begin{cases} 0 & \text{if } |x| \leq \lambda \\ -\infty & \text{otherwise} \end{cases}$. $\alpha < \lambda$, then 0, because $x = 0$ or $x > 0$. $\alpha \leq -\lambda$, $-\infty$, because $x = +\infty$. $\frac{\partial D}{\partial \alpha} = -\alpha - u \Rightarrow \alpha = \begin{cases} -u, & |u| \leq \lambda \\ -\lambda, & u \geq \lambda \\ \lambda, & u \leq -\lambda \end{cases}$. So $D(\alpha) = \begin{cases} -\frac{1}{2} \alpha^2 - \alpha u & \text{if } |\alpha| \leq \lambda \\ -\infty & \text{otherwise} \end{cases}$ max $D(\alpha)$. $\Rightarrow x = z = u + \alpha \Rightarrow$ so $x = S_\lambda(u) = \begin{cases} u - \lambda, & u \geq \lambda \\ 0, & \text{if } |u| \leq \lambda \\ u + \lambda, & u \leq -\lambda \end{cases}$.

6. If A p.s.d, then $\|Ax - b\|_2^2$ convex. 7. If $G(x)$ is convex, F is matrix, then $G(Fx)$ convex. 14. P2: $\min \frac{1}{2} \|\vec{x} - \vec{u}\|_2^2 + \lambda \|\vec{x}\|_1$, x_i are independent. $= \sum_{i=1}^K \left[\frac{1}{2} (x_i - u_i)^2 + \lambda |x_i| \right]$, so answer \vec{y} , $y_i = S_\lambda(u_i)$.

8. $\frac{\partial a^T x}{\partial x} = \frac{\partial x^T a}{\partial x} = a$. $(e^x)' = e^x$. $(\log x)' = \frac{1}{x}$. 15. P3. $\arg \min_x |x - u| + \lambda |x| \Rightarrow D(\alpha) = \min_{x, z} |z - u| + \lambda |z| + \alpha (x - z)$. $D(\alpha) = \min |z - u| + \alpha z + \min \lambda |z| + \alpha x$. $= \begin{cases} -\alpha u & \text{if } |\alpha| \leq 1 \\ -\infty & \text{otherwise} \end{cases} + \begin{cases} 0 & \text{if } |z| \leq \lambda \\ -\infty & \text{otherwise} \end{cases}$. \Rightarrow Dual Solution $\begin{cases} \alpha = -\min(1, \lambda), & u > 0 \\ \alpha = \min(1, \lambda), & u < 0 \end{cases}$. Suppose $u > 0$. Case 1. $\lambda > 1$, so $\alpha = -1 > -\lambda$, $x = 0$, from x sol. Case 2. $\lambda < 1$, so $\alpha = -\lambda$, $x = z = u$, from z sol. Case 3. $\lambda = 1$, so $\alpha = -1 = -\lambda$, $x = z \leq u$ from z sol, $x \geq 0$ from x sol, $0 \leq x \leq u$.

9. $\frac{\partial Ax}{\partial x} = A^T$. $\frac{\partial \|Ax - b\|_2^2}{\partial x} = \frac{\partial (Ax - b)^T (Ax - b)}{\partial x} = 2A^T(Ax - b)$. $\frac{\partial (x - s)^T A(x - s)}{\partial x} = (A + A^T)(x - s)$. $\phi(z) = 1/(1 + e^{-z})$. $\frac{\partial \phi}{\partial z} = \phi(z)/(1 - \phi(z))$. $(f \cdot g)' = f' \cdot g + g' \cdot f$. $(\frac{f}{g})' = \frac{f' \cdot g - g' \cdot f}{g^2}$. unconstrained optimization, check what happened as some components go to $\pm \infty$, check bounding condition if some components are bounded.

1. P4. $\arg\min_{\vec{x}} \frac{1}{2} \|\vec{x} - \vec{u}\|_2^2 + \lambda \|\vec{x}\|_1$
 $\arg\min_{\vec{x}} \frac{1}{2} \|\vec{x} - \vec{u}\|_2^2 + \lambda \|\vec{x}\|_1 + \vec{\alpha}^T (\vec{x} - \vec{z})$
 $\min_{\vec{x}} \frac{1}{2} \|\vec{x} - \vec{u}\|_2^2 - \vec{\alpha}^T \vec{z} + \min_{\vec{x}} \lambda \|\vec{x}\|_1 + \vec{\alpha}^T \vec{x}$
 $\vec{z} = \vec{u} + \vec{\alpha}$

$\min \lambda \|\vec{x}\|_1 + \vec{\alpha}^T \vec{x}$
 Change of vars. $\vec{x} = r\vec{y}$, $r = \|\vec{x}\|_1$. \vec{y} is vector with $\|\vec{y}\|_1 = 1$, so we min $\lambda r + r\vec{\alpha}^T \vec{y}$.
 Choosing a good \vec{y} , $\sum_{i=1}^K r \alpha_i y_i$, s.t. $|y_i| \leq 1$ for all i
 so $y = -1$ if $\alpha_i \geq 0 \Rightarrow \vec{y} = -\text{sign}(\vec{\alpha})$,
 +1 otherwise
 so $\sum r \alpha_i y_i = -r \|\vec{\alpha}\|_1 \Rightarrow \min \lambda r - r \|\vec{\alpha}\|_1$
 solution = $\begin{cases} 0 & \text{if } \lambda \geq \|\vec{\alpha}\|_1 \\ -\infty & \text{otherwise} \end{cases} \Rightarrow D(\vec{\alpha}) = \frac{1}{2} \|\vec{\alpha}\|_1^2 - \vec{\alpha}^T (\vec{u} + \vec{\alpha})$ s.t. $\|\vec{\alpha}\|_1 \leq \lambda$

2. P5. $\arg\min_{\vec{x}} \frac{1}{2} \|\vec{x} - \vec{u}\|_2^2 + \lambda \|\vec{x}\|_1$
 $\min_{\vec{x}} \frac{1}{2} \|\vec{x} - \vec{u}\|_2^2 + \lambda \|\vec{x}\|_1 + \vec{\alpha}^T (\vec{x} - \vec{z})$
 $\min_{\vec{x}} \frac{1}{2} \|\vec{x} - \vec{u}\|_2^2 - \vec{\alpha}^T \vec{z} + \min_{\vec{x}} \lambda \|\vec{x}\|_1 + \vec{\alpha}^T \vec{x}$
 $\vec{z} = \vec{u} + \vec{\alpha}$
 $r = \|\vec{x}\|_1$. \vec{y} unit vector, $\|\vec{y}\|_1 = 1$.
 $\min \lambda r + r\vec{\alpha}^T \vec{y}$, $r \geq 0$, $\|\vec{y}\|_1 = 1$.
 \vec{y} is in opposite direction of $\vec{\alpha}$. $y = -\frac{\alpha}{\|\alpha\|_1}$.
 so. $\min \lambda r - r \|\alpha\|_1$. so $r = \begin{cases} 0 & \text{if } \|\alpha\|_1 \leq \lambda \\ \infty & \text{otherwise} \end{cases}$
 $\Rightarrow \text{sol} = \begin{cases} 0 & \text{if } \|\alpha\|_1 \leq \lambda \\ -\infty & \text{otherwise} \end{cases}$

so $D(\vec{\alpha}) = -\frac{1}{2} \|\vec{\alpha}\|_1^2 - \vec{\alpha}^T \vec{u}$ s.t. $\|\vec{\alpha}\|_1 \leq \lambda$.
 $\vec{\alpha} = r\vec{y} \Rightarrow \max_{r,y} -\frac{1}{2} r^2 - r\vec{y}^T \vec{u}$ s.t. $r \leq \lambda$.
 $\vec{y} = -\frac{\vec{u}}{\|\vec{u}\|_1} \Rightarrow \max_{0 \leq r \leq \lambda} -\frac{1}{2} r^2 + r \|\vec{u}\|_1$
 $\frac{\partial}{\partial r} = -r + \|\vec{u}\|_1$. $r = \|\vec{u}\|_1 \Rightarrow r = \begin{cases} \|\vec{u}\|_1 & \text{if } \|\vec{u}\|_1 \leq \lambda \\ \lambda & \text{otherwise} \end{cases}$
 so $\vec{\alpha} = r\vec{y} = \begin{cases} -\vec{u} & \text{if } \|\vec{u}\|_1 \leq \lambda \\ -\lambda \vec{u} / \|\vec{u}\|_1 & \text{otherwise} \end{cases}$
 $\vec{x} = \vec{z} = \begin{cases} 0 & \text{if } \|\vec{u}\|_1 \leq \lambda \\ (1 - \lambda / \|\vec{u}\|_1) \vec{u} & \text{otherwise} \end{cases}$ switch that turn off x components

Partition x into groups: $[x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9]$
 $\arg\min \frac{1}{2} \sum \|x_{Gi} - u_{Gi}\|_2^2 + \sum \lambda_i \|x_{Gi}\|_1$
 soln = $x_{Gi} = \begin{cases} 0 & \text{if } \|u_{Gi}\|_1 \leq \lambda_i \\ (1 - \lambda_i / \|u_{Gi}\|_1) u_{Gi} & \text{otherwise} \end{cases}$

3. ADMM $\min f(x) + g(z) + \frac{\rho}{2} \|Ax + Bz - C\|_2^2$ s.t. $Ax + Bz - C = 0$
 $x_{t+1} \leftarrow \arg\min f(x) + \frac{\rho}{2} \|Ax + Bz_t - C\|_2^2 + \alpha_t^T Ax$ KKT: $Ax^* + Bz^* - C = 0$
 $z_{t+1} \leftarrow \arg\min g(z) + \frac{\rho}{2} \|Ax_{t+1} + Bz - C\|_2^2 + \alpha_t^T Bz$ $\vec{0} \in \partial_x f(x^*) + A^T \alpha^*$
 $\alpha_{t+1} \leftarrow \alpha_t + \rho(Ax_{t+1} + Bz_{t+1} - C)$ $\vec{0} \in \partial_z g(z^*) + B^T \alpha^*$

Consider x update:
 $0 \in \partial_x f(x_{t+1}) + \rho A^T (Ax_{t+1} + Bz_t - C) + A^T \alpha_t$ monitor convergence $\rho A^T B(z_{t+1} - z_t) \rightarrow 0$
 $0 \in \partial_x f(x_{t+1}) + A^T (\rho(Ax_{t+1} + Bz_t - C) + \alpha_t)$ Small ρ better?
 replace z_t with z_{t+1} , this is $\alpha_{t+1} - \rho Bz_{t+1} + \rho Bz_t$:
 $0 \in \partial_x f(x_{t+1}) + A^T (\alpha_{t+1} - \rho B(z_{t+1} - z_t))$ $\partial Ax_{t+1} + Bz_{t+1} - C = 0$
 $0 \in \partial_x f(x_{t+1}) + A^T \alpha_{t+1} - \rho A^T B(z_{t+1} - z_t)$ large ρ better
 $u = (\rho P)^{-1} \cdot \alpha$. $x_{t+1} = \arg\min (f(x) + \frac{\rho}{2} \|Ax + Bz_t - C + u\|_2^2)$
 $z_{t+1} = \arg\min (g(z) + \frac{\rho}{2} \|Ax_{t+1} + Bz - C + u\|_2^2)$, $u_{t+1} = u_t + (Ax_{t+1} + Bz_{t+1} - C)$

$f(x) = I_C(x)$. $x^* = \Pi_C(u)$. $f(x) = \lambda \|x\|_1$. $x^* = S_{\lambda/\rho}(u)$. $x^* = \begin{cases} u - \lambda/\rho \text{ sign}(u) & |u| > \lambda/\rho \\ 0 & \text{if } |u| \leq \lambda/\rho \end{cases}$

when $x^* = \arg\min f(x) + \frac{\rho}{2} \|x - u\|_2^2$
 4. $\min \|x\|_1$ subject to $Ax = b \Rightarrow \min \|x\|_1 + I_C(z)$. $C = \{x \mid Ax = b\}$, s.t. $x = z$
 $\min \|x\|_1 + I_C(z) + \frac{\rho}{2} \|x - z\|_2^2$. $u = \frac{1}{\rho} \cdot y$
 $\rho u^T (x - z) + \frac{\rho}{2} \|x - z\|_2^2 = \frac{\rho}{2} (\|x - z\|_2^2 + 2u^T (x - z) + \|u\|_2^2) - \frac{1}{2} \|u\|_2^2$
 $= \frac{\rho}{2} \|x - z + u\|_2^2 - \frac{\rho}{2} \|u\|_2^2$ so: $\min \|x\|_1 + I_C(z) + \frac{\rho}{2} \|x - z + u\|_2^2$

$x^{k+1} = \arg\min \|x\|_1 + \frac{\rho}{2} \|x - z^k + u^k\|_2^2 = S_{\lambda/\rho}(z^k - u^k)$
 $z^{k+1} = \arg\min I_C(z) + \frac{\rho}{2} \|x^{k+1} - z + u^k\|_2^2 = \Pi_C(x^{k+1} + u^k)$
 $u^{k+1} = u^k + (x^{k+1} - z^{k+1})$
 5. $\frac{\rho}{2} \|x - y + u\|_2^2 + \frac{\rho}{2} \|x - z + v\|_2^2 = \rho \|x - \frac{y+u}{2} - \frac{z+v}{2}\|_2^2$. $M = y - u$, $y = z - v$. $f(x) = \frac{1}{2} \|x - u\|_2^2 + \lambda \|x\|_1$

6. LASSO. $\min \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1$. $Fx - z = 0$.
 $L_P = \frac{1}{2} \|Ax - b\|_2^2 + \lambda \|x\|_1 + \frac{\rho}{2} \|Fx - z + u\|_2^2$
 $x^{k+1} = (A^T A + \rho F^T F)^{-1} (A^T b + \rho F^T (z^k - u^k))$
 $z^{k+1} = S_{\lambda/\rho}(F x^{k+1} + u^k)$: $u^{k+1} = u^k + F x^{k+1} - z^{k+1}$

7. Subgradient of f at \vec{z} is any vector \vec{v} s.t.
 $f(\vec{y}) \geq f(\vec{z}) + (\vec{y} - \vec{z})^T \cdot \vec{v}$ for all \vec{y} .
 If f is convex and $\vec{0} \in \partial f(x)$ then x is minimizer of f

8. SVM. $\min_w \lambda \sum_{i=1}^N (1 - y_i (w^T x_i + b)) + \frac{1}{2} \|w\|_2^2$, $A = \begin{bmatrix} -y_1 x_1 & \dots & -y_1 \\ \vdots & \ddots & \vdots \\ -y_m x_m & \dots & -y_m \end{bmatrix}$, $u = \begin{bmatrix} u \\ b \end{bmatrix} \in \mathbb{R}^{d+1}$
 $\Rightarrow \min \lambda \sum 1^T (A_i V_i + 1) + \frac{1}{2} \|z\|_2^2$ s.t. $V_i - z = 0$
 $\sum_{i=1}^B 1^T (A_i V_i + 1) + \frac{1}{2\lambda} \|z\|_2^2 + \sum_{i=1}^B \frac{\rho}{2} \|V_i - z + u\|_2^2$
 $u_i^{k+1} = \arg\min 1^T (A_i V_i + 1) + \frac{\rho}{2} \|V_i - z^k + u\|_2^2$
 $z^{k+1} = \arg\min \frac{1}{2\lambda} \|z\|_2^2 + \frac{\rho}{2} \sum \|u_i^{k+1} - z + u\|_2^2$
 $u_i^{k+1} = u_i^k + (V_i^{k+1} - z^{k+1})$
 s.t. $y_i (w^T x_i + b) \geq 1 - \xi_i$
 $\forall i, \xi_i \geq 0$

Variables:

1. x : observed variables (like data)
2. z : missing stuff / random variables we do not observe
3. θ parameters (not known and not random)

Model must specify $p(z|\theta)$, $p(x|z, \theta)$.

θ is usually split into 2 parts $p(z|\theta_1)$, $p(x|z, \theta_2)$

Goal: 1. Find value of θ so that model fits the data.

2. Find (approximate) posterior dist $p(z|x, \theta)$

both by maximizing likelihood $= p(x|\theta_1, \theta_2) = \sum_z p(x|z, \theta_2) p(z|\theta_1)$

$$\Lambda(q, \theta) = \sum_z q(z) \cdot \log \frac{p(z, x|\theta)}{q(z)}$$

if we could fully optimize q , then q would be $p(z|x, \theta)$.

because Λ is maximized when $KL(q||p(z|x, \theta))$ is minimized.

2. If we can't compute $p(z|x, \theta)$, then $q(z)$ is an approximation $p(z|x, \theta)$.

or if $z = (z_1, z_2)$, then $q(z) = q_1(z_1) q_2(z_2)$.

Exponential dist. $y \sim f(y|\beta) = \beta e^{-\beta y} \mathbb{1}_{\{y \geq 0\}}$

$$\int u dv = uv - \int u dv$$

$$\text{mean} = 1/\beta$$

$$\text{variance} = 1/\beta^2$$

$$= \int_0^\infty \beta y e^{-\beta y} dy$$

$$= -y e^{-\beta y} \Big|_0^\infty - \int_0^\infty e^{-\beta y} dy$$

$y \sim \text{Exponential}(\beta)$. z_1, \dots, z_n independently from $\text{Gauss}(y, \sigma^2)$.

x observed: x_1, \dots, x_n | z missing: y | θ param: β, σ^2

In this example, missing variable is only sampled once.

Full data likelihood $p(x, z|\theta) = p(x|z, \theta) p(z|\theta)$

$$= \left(\prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x_i - y)^2}{2\sigma^2}} \right) \beta e^{-\beta y} \mathbb{1}_{\{y \geq 0\}} \quad (*)$$

$\Lambda(q, \theta)$ function is $\sum_z q(z) \log \frac{p(x, z|\theta)}{q(z)}$

Set $q(y) = \text{Exponential}(\alpha) = \alpha e^{-\alpha y} \mathbb{1}_{\{y \geq 0\}}$

$$\Lambda(q, \theta) = \Lambda(q, \beta, \sigma^2) = \sum_z q(z) \log \frac{p(x, z|\theta)}{q(z)}$$

$$= \int_0^\infty \alpha e^{-\alpha y} \log \frac{(*)}{\alpha e^{-\alpha y}} dy$$

$$= \int_0^\infty \alpha e^{-\alpha y} \left[\sum_{i=1}^n \log \left(\frac{1}{\sqrt{2\pi}\sigma^2} \right) - \frac{(x_i - y)^2}{2\sigma^2} + \log \beta - y\beta - \log \alpha + \alpha y \right] dy$$

$$= -\frac{n}{2} \log(2\pi\sigma^2) + \log \beta - \log \alpha$$

$$+ \int_0^\infty \left[\sum_{i=1}^n \frac{1}{2\sigma^2} (x_i - y)^2 - y\beta + \alpha y \right] \alpha e^{-\alpha y} dy$$

$$\int_0^\infty y \alpha e^{-\alpha y} dy = \text{mean of exponential} = \frac{1}{\alpha}$$

$$\int_0^\infty (y - \frac{1}{\alpha})^2 \alpha e^{-\alpha y} dy = E(y - \frac{1}{\alpha})^2 = E(y - \frac{1}{\alpha} + \frac{1}{\alpha} - \frac{1}{\alpha})^2$$

$$= E(y - \frac{1}{\alpha})^2 + 2E(y - \frac{1}{\alpha}) \left(\frac{1}{\alpha} - \frac{1}{\alpha} \right) + \left(\frac{1}{\alpha} - \frac{1}{\alpha} \right)^2$$

$$\text{var} = 1/\alpha^2$$

$$= 1/\alpha^2 + (\frac{1}{\alpha} - \frac{1}{\alpha})^2$$

$$\text{So } \Lambda(\alpha, \beta, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) + \log \beta - \log \alpha - \frac{\beta}{\alpha} + 1 - \frac{1}{2\sigma^2} \sum_{i=1}^n \left[\frac{1}{\alpha^2} + (x_i - \frac{1}{\alpha})^2 \right]$$

Project onto a convex set C : $\argmin_x \frac{1}{2} \|x - y\|_2^2$ s.t. $Ax = b$

project onto set $\{x | Ax = b\}$. $\Lambda = \frac{1}{2} \|x - y\|_2^2 + \alpha^T (Ax - b)$

dual problem $\argmax_\alpha [\argmin_x \Lambda(x, y)]$

$$\argmin_x \frac{1}{2} \|x - y\|_2^2 + \alpha^T (Ax - b)$$

$$\frac{\partial}{\partial x} = x - y + A^T \alpha = 0 \Rightarrow x = y - A^T \alpha$$

$$\text{plug in: } \frac{1}{2} \|A^T \alpha\|_2^2 + \alpha^T A y - \|A^T \alpha\|_2^2 - \alpha^T b$$

$$= -\frac{1}{2} \|A^T \alpha\|_2^2 + \alpha^T A y - \alpha^T b, \quad \frac{\partial}{\partial \alpha} = -A A^T \alpha + A y - b = 0$$

$$\text{Since } x = y - A^T \alpha = y - A^T (A A^T)^{-1} (A y - b)$$

special case A is a vector w^T , $x = y - \frac{b}{\|w\|_2^2} (y \cdot w - b)$

if $w \cdot y = b$ then $x = y$

$$(*) \argmin_x \frac{1}{2} \|x - y\|_2^2 \text{ s.t. } w \cdot x \leq b$$

if $w \cdot y > b$, then x is a boundary

(boundary is convex) $w \cdot x = b$

$$\argmin_x \frac{1}{2} \|x - y\|_2^2 \text{ s.t. } w_1 x \leq b_1, w_2 x \leq b_2, w_3 x \leq b_3 \dots \text{ now use ADMM.}$$

$$\argmin_x \frac{1}{2} \|x - y\|_2^2 + I_{\{w_1 x \leq b_1\}} + I_{\{w_2 x \leq b_2\}} + \dots \text{ s.t. } x = z$$

$$x \leftarrow \argmin_x \frac{1}{2} \|x - y\|_2^2 + \frac{\rho}{2} \|x - z + \alpha\|_2^2 + \frac{\rho}{2} \|x - 0 + \alpha\|_2^2 + \dots$$

$$z \leftarrow \argmin_z I(w, z \leq b) + \frac{\rho}{2} \|x - z + \alpha\|_2^2 \quad \text{projection of } x + \alpha \text{ onto } (*)$$

$$p \sim \text{Beta}(a, b) \text{ density } f(p) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} p^{a-1} (1-p)^{b-1} \mathbb{1}_{\{p \in [0, 1]\}}$$

$$\text{mean} = \frac{a}{a+b} \quad \text{variance} = \frac{ab}{(a+b)^2 (a+b+1)}$$

$$p \sim \text{Beta}(a, b), x \sim \text{Bernoulli}(p)$$

$$\text{Joint: } p(x=u, p=p) = \binom{n}{k} p^k (1-p)^{n-k} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} p^{a-1} (1-p)^{b-1}$$

$$= \binom{n}{k} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(k+a)\Gamma(n-k+b)}{\Gamma(n+a+b)} \quad p(p|k) = \frac{p(p, k)}{p(k)} = \frac{p^{k+a} (1-p)^{n-k+b}}{p^k (1-p)^{n-k}} = \frac{p^{a+b} (1-p)^{n-k+b-k}}{p^{a+b} (1-p)^{n-k+b-k}} = \text{Beta}(k+a, n-k+b)$$

$$f(p) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} p^{a-1} (1-p)^{b-1}$$

$$\int_0^1 f(p) dp = 1 \Rightarrow \int_0^1 p^{a-1} (1-p)^{b-1} dp = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

conjugate: $p(p|k)$ is the same dist as p | $KL(q||p) = \sum_z q(z) \log \frac{q(z)}{p(z)}$

$$KL(q||p) \geq 0, KL(q||q) = 0, \argmin_q KL(q||p) = p, \argmin_p KL(q||p) = q$$

- #12. Fixed cluster centers μ_1, \dots, μ_k (unknown)
 to generate a point x_j , ①. pick a cluster i with prob π_i (multinomial $(1, \pi)$). ②. $x_j \sim \text{Gauss}(\mu_i, 1)$
 1. observed: x_1, \dots, x_n . 2. missing: cluster of x_j , call it z_j .
 3. unknown: $\Pi(\pi_1, \dots, \pi_k), \mu_1, \dots, \mu_k$.
 $z_j = [0, \dots, 1, \dots, 0]$ 1 word is 1, everything else is 0.
 $z_j[i] = 1$ means x_j is in cluster i .

complete data likelihood function
 $p(x_1, z_1, x_2, z_2, \dots, x_n, z_n | \mu_1, \dots, \mu_k, \pi)$
 $= \prod_{j=1}^n p(x_j | z_j) p(z_j | \mu_1, \dots, \mu_k, \pi)$
 if cluster of point j is 1 then $\frac{1}{\sqrt{2\pi}} e^{-\frac{(x_j - \mu_1)^2}{2}} \pi_1$
 if cluster of point j is 2 then $\frac{1}{\sqrt{2\pi}} e^{-\frac{(x_j - \mu_2)^2}{2}} \pi_2$
 $= \prod_{j=1}^n \prod_{i=1}^k \left[\frac{1}{\sqrt{2\pi}} e^{-\frac{(x_j - \mu_i)^2}{2}} \pi_i \right]^{z_j[i]}$

Full data loglikelihood
 $\sum_{j=1}^n \sum_{i=1}^k z_j[i] \left[-\frac{(x_j - \mu_i)^2}{2} + \log \pi_i \right] + \text{const.}$
 $q(z_1, \dots, z_n) = q_1(z_1) q_2(z_2) \dots q_n(z_n)$

$\Lambda(q, \mu_1, \dots, \mu_k, \pi) =$
 $\sum_{z_1, \dots, z_n} q_1(z_1) \dots q_n(z_n) \sum_{j=1}^n \sum_{i=1}^k z_j[i] \left[-\frac{(x_j - \mu_i)^2}{2} + \log \pi_i \right]$
 $= \sum_{z_1, \dots, z_n} q_1(z_1) \dots q_n(z_n) \left[\log q_1(z_1) - q_n(z_n) \right]$

Simplifying Steps:
 ① $\sum_{z_1, \dots, z_n} q_1(z_1) \dots q_n(z_n) \left[\log q_1(z_1) - q_n(z_n) \right]$
 $= \sum_{z_1, \dots, z_n} q_1(z_1) \dots q_n(z_n) \log q_1(z_1) + \dots$
 $= \sum_{z_1, \dots, z_n} q_1(z_1) q_2(z_2) \dots q_n(z_n) \log q_1(z_1)$
 $= \sum_{z_1} q_1(z_1) \log q_1(z_1) \left[\sum_{z_2, \dots, z_n} q_2(z_2) \dots q_n(z_n) \right]$
 $= \sum_{z_1} q_1(z_1) \log q_1(z_1) \cdot 1$
 so ① = $\sum_{z_1} q_1(z_1) \log q_1(z_1) + \sum_{z_2} q_2(z_2) \log q_2(z_2) + \dots$

note here z_i has k numbers
 multinomial (n, p) vector $p(k) = \left(\frac{n!}{\prod_i k_i!} \right) \prod_i p_i^{k_i}$
 vector of counts, mean = np^e

$\sum_{z_1, \dots, z_n} q_1(z_1) \dots q_n(z_n) \sum_{j=1}^n \sum_{i=1}^k z_j[i] \left[-\frac{(x_j - \mu_i)^2}{2} + \log \pi_i \right]$
 $= \sum_{z_1, \dots, z_n} q_1(z_1) \dots q_n(z_n) \sum_{i=1}^k z_i[i] \left[-\frac{(x_i - \mu_i)^2}{2} + \log \pi_i \right] + \dots$
 $= \sum_{z_1} q_1(z_1) \sum_{i=1}^k z_i[i] \left[-\frac{(x_i - \mu_i)^2}{2} + \log \pi_i \right], \text{ so}$
 $\Lambda = \sum_{z_1} q_1(z_1) \sum_{i=1}^k z_i[i] \left[-\frac{(x_i - \mu_i)^2}{2} + \log \pi_i \right] - \sum_{z_1} q_1(z_1) \log q_1(z_1) +$
 $+ \sum_{z_2} q_2(z_2) \sum_{i=1}^k z_i[i] \left[-\frac{(x_i - \mu_i)^2}{2} + \log \pi_i \right] - \sum_{z_2} q_2(z_2) \log q_2(z_2) + \dots$
 $= h(z)$ $\leftarrow q$ update.

by KL-divergence trick, $q_i(z_i) = \frac{e^{h(z_i)}}{\sum_r e^{h(z_i)}}$
 probability of cluster j : $q_i(z_i = [0, \dots, 1, \dots, 0]) = \frac{\pi_j e^{-\frac{(x_i - \mu_j)^2}{2}}}{\sum_{s=1}^k \pi_s e^{-\frac{(x_i - \mu_s)^2}{2}}} = T_j[i]$
 $z_i[i]$ is 0 or 1, when $z_i = [1, 0, \dots, 0]$ is 1 = $T_1[i]$, prob $q_i(z_i = [1, 0, \dots, 0])$
 plugin: $\sum_{z_1} q_1(z_1) \sum_{i=1}^k z_i[i] \left[-\frac{(x_i - \mu_i)^2}{2} + \log \pi_i \right]$
 $= \sum_{i=1}^k \left[-\frac{(x_i - \mu_i)^2}{2} + \log \pi_i \right] \sum_{z_1} q_1(z_1) z_i[i]$
 $= \left[-\frac{(x_1 - \mu_1)^2}{2} + \log \pi_1 \right] \sum_{z_1} q_1(z_1) z_1[1] + \dots$

so: $\sum_{z_1} q_1(z_1) \sum_{i=1}^k z_i[i] \left[-\frac{(x_i - \mu_i)^2}{2} + \log \pi_i \right] = \sum_{i=1}^k T_i[i] \left[-\frac{(x_i - \mu_i)^2}{2} + \log \pi_i \right]$
 after q step, $\Lambda = \sum_{i=1}^k T_i[i] \left[-\frac{(x_i - \mu_i)^2}{2} + \log \pi_i \right] + \text{something with no } \mu \text{ or } \pi$
 $+ \sum_{i=1}^k T_i[i] \left[-\frac{(x_i - \mu_i)^2}{2} + \log \pi_i \right] + \text{something}$

μ_i update: $\mu_i = \text{argmax}_{\mu_i} \sum_{j=1}^n T_j[i] \left[-\frac{(x_j - \mu_i)^2}{2} \right]$
 $\frac{\partial}{\partial \mu_i} = \sum_{j=1}^n T_j[i] (x_j - \mu_i) = 0$
 π update: $\pi = \text{argmax}_{\pi} \sum_{j=1}^n \sum_{i=1}^k T_j[i] \log \pi_i$ s.t. $\sum_{i=1}^k \pi_i = 1$, Lagrange multiplier λ .
 $\frac{\partial \pi_i}{\partial \pi_i} = \sum_{j=1}^n T_j[i] / \pi_i + \lambda = 0$
 $\Rightarrow \pi_i \lambda + \sum_{j=1}^n T_j[i] = 0$
 $\Rightarrow \pi_2 \lambda + \sum_{j=1}^n T_j[2] = 0$, add them up using $\sum \pi_i = 1$.
 to get $\lambda + \sum_{j=1}^n \sum_{i=1}^k T_j[i] = 0$, so $\lambda + n = 0$, $\lambda = -n$.
 so $\pi_i = \frac{\sum_{j=1}^n T_j[i]}{n}$

$\Sigma = \text{Diag}(np) - np p^T$
 $E[(x[i] - \mu[i])(x[j] - \mu[j])] = -np[i] p[j]$
 vector, $E[(x[i] - \mu[i])^2] = np[i](1 - p[i])$