Note: Due to the nature of the tutorials, things are a bit disjointed.
I have labeled most of the stuff.
Good luck everyone!   Except Issam, you don't need luck =)        — Scott

## Some Matrix Properties

$$f(x) = x^T A x + bx + c$$

$$\nabla f(x) = Ax + b$$

ie: $\frac{1}{2}(A + A^T) + b$

symetric

$$f(x) = \| Ax - b \|^2$$

$$= \frac{1}{2}\left( b^T b - 2x^T A^T b + x^T A^T A x \right)$$

$$\nabla f(x) = 0 - A^T b + A^T A x$$

$$= A^T (Ax - b)$$

$$\sum_i v_i x_i = v^T x = x^T v$$

$$v = A^T B C Y_w$$

$$(n \cdots x_1)$$

$$\sum_i x_i \sum_j x_j a_{ij} = x^T A x$$

$$\sum_i x_i v_i = x^T v$$

$$v_j = x_j a_{ij}$$

$$A x = \begin{bmatrix} \sum_j x_j a_{1j} \\ \vdots \\ \sum_j x_j a_{nj} \end{bmatrix}$$

$$f = x^T A x$$

$$\nabla f = A x$$

$$\nabla^2 f = A$$

$$\min_w f(w) + \frac{\lambda}{2}\|w\|^2$$

$$w^{t+1} = w^t - \alpha\left( \nabla f(w^t) + \lambda w^t \right)$$

## Grad descent

$$w^{t+1} = w^t - \alpha \nabla g(w^t)$$

$$\underset{d \text{ by } 1}{w^t} = \underset{1 \text{ by } 1}{\beta^t}\, \underset{d \text{ by } 1}{v^t}$$

ex

$$w^t - \alpha \lambda w^t - \alpha \nabla f_i(w^t) = (1 - \alpha\lambda) w^t - \alpha f_i(x^t)$$

$$\underset{1 \text{ by } 1}{} \quad \underset{d \text{ by } 1}{}$$

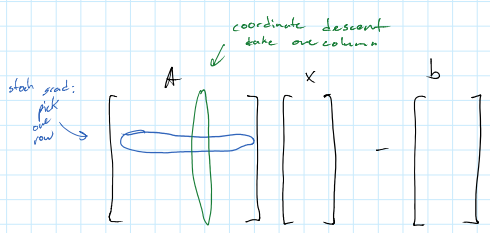$$w^{t+1} = \overset{\text{proportional}}{(1 - \alpha\lambda)} w^t \qquad O(d)$$

$$\beta^{t+1} v^{t+1} = (1 - \alpha\lambda)\beta^t v^t$$

$$= \left[ (1 - \alpha\lambda)\beta^t \right] v^t$$

$$v^{t+1} = v^t$$

$$\beta^{t+1} = (1 - \alpha\lambda)\beta^t \qquad O(1)$$

## Stoch Grad vs Coord Desc.

$V \geqslant \max\{x_i\} \Rightarrow V \geqslant x_i \quad \forall i$

$V \geqslant \max(a, b)$

$\quad \Rightarrow V \geqslant a,$
$\qquad V \geqslant b$

$$\min_{w, V, r} \sum_{i=1}^{N} V_i + \lambda r \qquad s.t.$$

$\qquad V_i \geqslant w^T x_i - y_i$
$\qquad V_i \geqslant y_i - w^T x_i$

$\qquad r \geqslant w_3$
$\qquad r \geqslant -w_3$

---

LP: $\displaystyle\min_x c^T x$
$\qquad\qquad s.t. \quad Ax \leqslant b$

QP: $\displaystyle\min_x \frac{1}{2} x^T H x + c^T x$
$\qquad\qquad\qquad s.t. \quad Ax \leqslant b$

---

prox: $\qquad \displaystyle\min_x f(x) + r(x) \qquad$ i.e. $r(x)$ non continuous

$$x^{t+1} = prox_{\alpha r(\cdot)} \left[ x^t - \alpha_t \nabla f(x^t) \right]$$

$$prox_{\alpha r(\cdot)} [y] = \arg\min_x \frac{1}{2} \|x - y\|^2 + \alpha r(y)$$

proj

$$r(x) = \begin{cases} 0 & if \quad Ax \leqslant b \\ \infty & if \quad Ax > b \end{cases}$$

QP: $\displaystyle\min_x \underbrace{\frac{1}{2} x^T H x + c^T x}_{f(x)} + r(x)$

---

## On time to error

$$\|x^t - x^*\| = O\left(\frac{1}{t}\right) \qquad \begin{array}{l}\text{how big } t, \\ \text{such that } x^t - x^* \leqslant \epsilon\end{array} \qquad \begin{array}{ll} O\left(\frac{1}{t}\right) & \leqslant \frac{1}{t} \\ O(p^t) & \leqslant c p^t \end{array}$$

$$\leqslant c \frac{1}{t}$$

$$\leqslant \epsilon$$

$\dfrac{c}{t} \leqslant \epsilon \quad \Rightarrow \quad \dfrac{c}{\epsilon} \leqslant t \qquad\qquad t \geqslant O\left(\frac{1}{\epsilon}\right)$

$$t = \Omega\left(\frac{1}{\epsilon}\right)$$

---

## Dual norms

$f(x) = \|x\|_p$

$$f^*(y) = \begin{cases} 0 & if \quad \|y\|_q \leqslant 1 \\ \infty & else \end{cases} \qquad \begin{array}{l} 1 \to \infty \\ 2 \to 2 \end{array} \qquad \frac{1}{p} + \frac{1}{q} = 1$$

$f(x) = \|x\|_p^2 \qquad\qquad \lambda \|x\|_p^2$

$f^*(y) = \|y\|_q^2 \qquad\qquad \frac{1}{\lambda} \|y\|_q^2$

if $\quad g(x) = a f(x)$

then $\qquad g^*(y) = a f\left(\frac{y}{a}\right)$

<u>conjugate</u>

$$f^*(y) = \sup_x \{ y^+ x - f(x) \}$$

$$f(x) = \exp(x)$$

$$f^*(y) = \sup_x \{ \underbrace{yx - \exp(x)}_{g(x)} \}$$

$$\nabla g(x) = y - \exp(x)$$

at minimum , $\nabla g(x) = 0$

$$0 = y - \exp(x)$$

solve for $x$

$$x = \log(y) \qquad \left( \text{for } y > 0 \right)$$

plug $x$ back in

$$f^*(y) = y \log(y) - \exp(\log(y))$$

$$= y \log y - y$$

$$f^*(y) = y(\log(y) - 1) \qquad \text{domain : } y > 0$$

$$f^*(y) = \begin{cases} y(\log(y) - 1) & y > 0 \\ \infty & y \leq 0 \end{cases}$$

---

<u>"separable" function</u>

$$\min \sum_{i=1}^{d} f_i(x_i) \qquad \text{different } f \text{ to each variable}$$

ex: $\quad f(x) = x_1 + x_2^2 + \exp(-x_3) + \tanh(x_4) \cdots$

$$\min_{x_1 x_2 x_3 x_4 \cdots} f(x) \qquad \text{is independent}$$

equivalent:          — can't have $x_1 = x_2$

$$\sum_{i=1}^{d} \min_{x_i} \{ f(x_i) \}$$

---

<u>EM:</u>



$$Q(\theta | \theta^t) = \sum_i \sum_h \text{with } \log p(x, h | \theta)$$

Tutorial Part 2

<u>EM:</u>   $Q(\theta|\theta^t) \gtrless \sum_{i=1}^{N} \sum_{h} w_{ih} \log p(x,h|\theta)$
(will be bonus)
$\uparrow \theta^t$

— sum outside log, not inside

<u>Prisoners dilema:</u>

$p(F|I) = 0.01$   footprints, given innocent

$p(\neg F|I) = 1 - p(F|I)$

$p(F|I) \propto p(I|F) p(F)$
$\hookrightarrow 1 - p(\neg I|F)$

<u>Algorithms</u>
— may ask about runtime

```
O O O
O O O         decision stumps?.

X X X
```

```
O O O    X
O O O   X
       X  X    how  about now?.
X X X   X X      — maybe with boosts
```

(interpret plot)

<u>techniques</u>                              ← maybe 1 norm?.
                                             some other norm
          1.   $f(x^{k+1}) \geq f(x^k) - \frac{1}{2\mu} \| \nabla f(x^k) \|^2$

likely on
test.     2.   $f(x^*) \leq f(x^k) - \frac{1}{2L} \| \nabla f(x^k) \|^2$

         combine inequalities for solution

         $f(x^{k+1}) - f(x^*) \leq \left(1 - \frac{\mu}{L}\right) \left[ f(x^k) - f(x^*) \right]$

         $\|x^{k+1} - x^*\|^2$
                    $= \| (x^k - \alpha_k \nabla f(x^k)) - x^* \|^2$
                    $= \| (x^k - x^*) - \alpha_k \nabla f(x^k) \|^2$

on test: replace $l_2$ with $l_1$

          ber with multi
          etc . . .          (small changes to questions you've seen)

coding:      implement...      ex    N B

KNN, with $l_1$ norm

start/end   of   code   similar   to   assn

i.e.   will   have   a   start   point

Define   stuff   on   your   cheat   sheet   if   you   don't   know   it!

ex: "what   is   naive   bayes"

## Lin Programing

$\frac{1}{2} x^T A x + b^T x$

$\frac{1}{2} ||Ax - b||^2$

$\sum_{i=1}^{N} z_i (w^T a_i - b)^2$

$\Bigg\}$ Quadratic Functions

Not Linear Constraint:   (bad)

$V_i \geqslant \max \{ x_{ij} - x_i \}$

Linear Constraint:   (good)

$V_i \geqslant x_i$

$V_i \geqslant -x_i$

$\min_x \sum_{i=1}^{d} |x_i|$

$\Downarrow$

non-smooth   to   smooth   conversion

$\Downarrow$

$\min_{x, v} \sum_{i=1}^{N} V_i$    subject to    $V_i \geqslant x_i$

$V_i \geqslant -x_i$

linear constraints! good!

Norm   equivalences:   You   can   use   them   directly
(so maybe write them down!)

Hessian:   outer product

$\nabla^2 f(x) = \sum_i a_i a_i^T$    "least squares Hessian"    least squares:  $f(x) = \frac{1}{2} ||Ax - b||^2$

$\underset{d \times 1}{} \quad \underset{1 \times d}{} \quad = d \times d$

$\begin{bmatrix} \sum_j a_{ij} a_{ij} \\ \\ \\ \end{bmatrix}$

$\Downarrow$

$\begin{bmatrix} a_1^T a_1 & a_1^T a_2 \\ a_2^T a_1 & a_2^T a_2 \end{bmatrix} = A^T A$

Lets say:

← 1×1 scalar

$$\nabla^2 f(x) = \sum_i a_i a_i^T d_i$$

$$= \sum_i (a_i d_i) a_i^T$$

$$\begin{bmatrix} \sum_j a_{1j} d_j a_{1j} & \sum_j a_{1j} d_j a_{2j} & \sum_j a_{1j} d_j a_{3j} \\ \sum_j a_{2j} d_j a_{1j} & \cdots & \\ \vdots & & \end{bmatrix}$$

$$\begin{bmatrix} a_{11} d_1 \\ a_{12} d_2 \\ a_{13} d_3 \\ \vdots \end{bmatrix} \qquad \begin{bmatrix} (a_1 \cdot d)^T a_1 & (a_1 \cdot d)^T a_2 \\ & \\ & \end{bmatrix}$$

$$= A^T \operatorname{diag}(d) A \qquad \begin{bmatrix} d_1 & & & 0 \\ & d_2 & & \\ & & \ddots & \\ 0 & & & \end{bmatrix} \qquad \begin{bmatrix} d_1(a_{11} \ a_{12} \cdots) \end{bmatrix}$$

a trick:

$$f(Ax)$$

$$\nabla f(Ax) = A^T \nabla f(x) \qquad \nabla f(x) = A^T \begin{bmatrix} \frac{\partial f}{\partial x_i} \end{bmatrix}$$

$$\nabla^2 f(Ax) = A^T \operatorname{diag}(\nabla^2 f(x)) A \qquad A^T \operatorname{diag}\left[ \frac{\partial^2 f}{\partial x_i} \Big|_{Ax} \right] A$$

$$\underbrace{\qquad}_{\sum_{i=1}^N f_i(a_i^T x)}$$

partial derivative evaluated at Ax

Logistic

$$f(x) = \sum_{i=1}^N \log\left(1 + \exp\left(-b_i a_i^T x\right)\right) = f(Ax)$$

$$\nabla f(x) = A^T \left[ y - \frac{y}{1 + \exp(b_i a_i^T x)} \right] \qquad \text{or something similar}$$

$$\qquad\qquad b\theta(w)$$

$$\nabla^2 f(x) = A^T \operatorname{diag}(\text{something}) A$$

$$\nwarrow \ \theta(w)(1 - \theta(w))$$

write/memorize operations that preserve convexity!
    - don't want to take hessians to proove convexity

another trick (probably not on midterm, but you never know)

$$\frac{1}{2} \| Ax - b \|^2 + \lambda \| x \|_1$$

$$\Downarrow$$

$$\min_{x^+, x^-} \frac{1}{2} \| A(x^+ - x^-) - b^2 + \lambda \sum_j (x_j^+ + x_j^-) \qquad \text{quadratic}$$

$$\text{s.t.} \quad x_j^+ \geq 0 \qquad\qquad \text{linear}$$
$$\qquad\qquad x_j^- \geq 0$$

at solution
$$x = x^+ - x^- \qquad\qquad x \in \mathbb{R}^d$$
$$\qquad\qquad\qquad x^+ \in \mathbb{R}^d$$
$$\qquad\qquad\qquad x^- \in \mathbb{R}^d$$

at solution
$$x^+ = x \, \mathcal{I}(x > 0)$$

$$x^- = x \, \mathcal{I}(x < 0)$$

derive co-ordinate descent

derive along a dimension ($x_j^+$, or $x_j^-$), set to zero, solve,
...
yadda yadda

linear constraints, same formulation as dual sum

dual to primal

$$f(Ax) + g(x)$$

$$-f^*(-y) - g(A^\top y)$$

if I now have optimal $y$ now? How to get $x$?

remember: $\sup_x \{ y^\top x - f(Ax) \}$

$x$ that solves this

No non-differentiable functions for the conj on midterm
(but maybe norms! or convert to smooth!)

$$f(x) = \|x\|_p$$
$$f^*(y) = \begin{cases} \infty & \text{iotw} \quad \|y\|_q \leq 1 \end{cases}$$

$$f(x) = \frac{\lambda}{2} \|x\|_p^2 \quad \Rightarrow \quad f^*(y) = \frac{1}{2\lambda} \|y\|_q^2$$

nLL                reg

Fair game question:

$$\log(p(w|y,x)) = \sum_{i=1}^{N} \log\left(1 + \exp\left(-b_i a^\top x\right)\right) + \frac{\lambda}{2}\|w\|^2$$

possible question:
show convex!

(log-likelihood)          (log-prior)
                          (regularizer)

$$\exp(a+b)$$
$$= \exp(a)\exp(b)$$

$$\exp\left(-\sum_{i=1}^{N} \log\left(1 + \exp\left(-b_i a^\top x\right)\right)\right)$$

$$\exp\left(-\frac{\lambda}{2}\|w\|^2\right) \quad \leftarrow \text{gaussian}$$

$$w_j \sim N(0, \lambda^{-1})$$

$$\prod_i^{N} \exp\left(-\log\left(1 + \exp\left(-b_i a^\top x\right)\right)\right)$$

others:
$\ell_1$ laplacian
regularizer could be logistic!

$$\prod_i^{N} \frac{1}{1 + \exp\left(-b_i a^\top x\right)}$$

logistic function

$$\sigma\left(b_i a_i^\top x\right) \quad \text{sigmoid}$$

should be able to go both ways!

Converting to smooth, could ask:

$$+ \lambda \|w\|_\infty$$

$$= \max_i \{ |w_i| \}$$

$\max_i \& \max \{ \cdots \}\}$

one variable that bounds
also bounds the absolute!

<u>Multinomial</u> could be on midterm!

$$\text{Mult}(x|\theta) = \theta_1^{\mathbb{I}[x=1]} \; \theta_2^{\mathbb{I}[x=2]} \; \theta_3^{\mathbb{I}[x=3]}$$

$$\text{Dir}(\theta|a) \propto \theta_1^{a_1-1} \theta_2^{a_2-1} \theta_3^{a_3-1}$$

$$p(\theta|x,a) \propto p(x|\theta) p(\theta|a)$$

$$= \theta_1^{\mathbb{I}[x=1]+a_1-1} \; \theta_2^{\mathbb{I}[x=2]+a_2-1} \; \theta_3^{\mathbb{I}[x=3]+a_3-1}$$

$$\theta_1^{(\mathbb{I}[x=1]+a_1)-1} \; \theta_2^{(\mathbb{I}[x=2]+a_2)-1} \; \theta_3^{(\mathbb{I}[x=3]+a_3)-1}$$

$$p(\theta|x,a) \sim \text{Dir}(\mathbb{I}[x=1]+a_1, \; \mathbb{I}[x=2]+a_2, \; \mathbb{I}[x=3]+a_3) \qquad \text{\color{blue}posterior, prior, same distribution family}$$

$$p(x|n,\theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x} \qquad\qquad \color{green}\propto \theta^x (1-\theta)^{n-x}$$

$$p(\theta|a,b) = \frac{1}{B(a,b)} \theta^{a-1}(1-\theta)^{b-1} \qquad \color{green}\propto \theta^{a-1}(1-\theta)^{b-1}$$

$$\color{green}p(\theta|x,n,a,b) \propto \theta^x(1-\theta)^{n-x} \theta^{a-1}(1-\theta)^{b-1}$$

$$\color{green}\Rightarrow \theta^{(x+a)-1}(1-\theta)^{(n-x+b)-1}$$

$$\color{green}\theta|x,n,a,b \sim \text{Beta}(x+a, \; n-x+b)$$

<u>Back to inf norm</u>

$$\min_x \; \tfrac{1}{2}\|Ax-b\|^2 + \lambda\|x\|_\infty$$

$\Downarrow$  1. write in terms of "max"

$$\min_x \; \tfrac{1}{2}\|Ax-b\|^2 + \lambda \max_j \{|x_j|\}$$

$$\min_x \; \tfrac{1}{2}\|Ax-b\|^2 + \lambda \max_j \{\max\{x_j, -x_j\}\}$$

$\Downarrow$  2. upper bound max by linear variable

$$\min_{x,v} \; \tfrac{1}{2}\|Ax-b\| + \lambda v \qquad \text{subject to} \quad v \geq \max_j\{\max\{x_j, -x_j\}\}$$

$\Downarrow$  3. $v \geq \max\{a,b\} \Rightarrow v \geq a, \; v \geq b$

$$\min_{x,v} \; \tfrac{1}{2}\|Ax-b\| + \lambda v \qquad \text{subject to}$$
$$v \geq x_j$$
$$v \geq -x_j$$

<u>Midterm Info</u>   all questions   equally   weighted,   but harder questions   later

likely:

<u>8 questions</u>

(bonus worth half a question)


difficulty / question

Example Midterm Questions:

show convex using convex properties      (medium difficulty)

bayes rule!      (easy difficulty)

converting non-smooth to smooth      (medium difficulty)

fenchel dual      (medium-hard difficulty)