

An example, with a variety of algorithms

$$\min_x \|y - Ax\| + \lambda \|x\|_1$$

\uparrow
 $\sum_i |x_i|$

This problem turns up often

- Sparse reconstruction

- y is a signal
- A is a dictionary
- x is a sparse encoding

l1 reg linear regression

model y_i by $\alpha_i^T x + \alpha_b$
the bias

solve

$$\sum_i (y_i - (\alpha_i^T x + \alpha_b))^2 + \lambda |x|_1$$

(it is not usual to regularize the bias.
- you could estimate it separately)

$$\|y - (Ax + \alpha_b)\|^2 + \lambda |x|_1$$

Now estimate α_b separately, to

get \tilde{y}

$$\|\tilde{y} - Ax\|^2 + \lambda |x|_1$$

In both cases, the attraction is that the L_1 norm encourages sparsity (many zeros in x). One way to see this is notice that the penalty for small x_i is large compared to L_2 — its worth making small values zero.

Compressed Sensing:

• assume we have m linear measurements of an unknown signal z

$$y_i = \phi_i \cdot z + \text{noise}$$

④

Suppose we know that z is compressible, or has a sparse repn in a transform domain ~~(with a dictionary)~~ W

then we can recover z

by solving

$$\| \Phi z - y \|^2 + \lambda \| Wz \|_1$$

↓ measurement values

↑ measurement vectors, forming the compressed sensing matrix

↑ sparsifying transform.

Now assume that W is invertible

5

we can solve

$$\min \|Ax - y\|^2 + \lambda \|x\|_1$$

$$\text{where } Wx = z$$

$$AW = \Phi$$

A is the dictionary

This problem is VERY DIFFERENT from

$$\|y - Ax\|^2 + \lambda \|x\|^2$$

↑ 2 norm.

2-norm problem:

$$(A^T A + \lambda I) x = A^T y$$

→ linear!

1-norm:

not linear

2 norm:

$$\lim_{\lambda \rightarrow 0} x = \frac{I^{-1} A^T y}{(AA)^{-1} A^T y}$$

$$= (A^T A)^+ A^T y$$

+ → Moore-Penrose pseudo inverse

2-norm: $\lambda \rightarrow \infty$ gives $x \rightarrow 0$

1-norm: ~~$\lambda \rightarrow \infty$ gives~~

$x = 0$ for values of $\lambda \geq \lambda_{\max}$

where $\lambda_{\max} = \|2A^T y\|_{\infty}$

($\|u\|_{\infty} = \max_i |u_i|$, inf norm).

2-norm

$x(\lambda)$ is a curve

(rational, algebraic, some

wildly interesting geometry)

1-norm

$x(\lambda)$ is piecewise linear!

Some transformations of this problem

$$\min \|Ax - y\|^2 + \lambda |x|_1$$

can be turned into a quadratic program.

$$\min \|Ax - y\|^2 + \lambda \sum_i t_i$$

st. $-t_i \leq x_i \leq t_i$


|||

$$-x_i - t_i \leq 0$$

$$x_i - t_i \leq 0$$

Notice this might be worrying - the

objective is $x^T A^T A x - 2y^T A x + y^T y$

This could be  sense

9

Consider

$$\min [f(x) + \lambda g(x)]$$

1

and

$$\min f(x)$$

2

$$\text{st } g(x) \leq \mu.$$

Lagrangian for (2):

$$f(x) + \lambda (g(x) - \mu)$$

KKT

$$\nabla f + \lambda \nabla g$$

$$\lambda \geq 0$$

$$g\lambda = 0$$

$$g(x) - \mu = 0$$

~~exclude the~~

Different values of $\mu \rightarrow$ values of $\lambda \geq 0$

So for any $\lambda \geq 0$ in ①

I can choose a $\mu \geq 0$ in ②

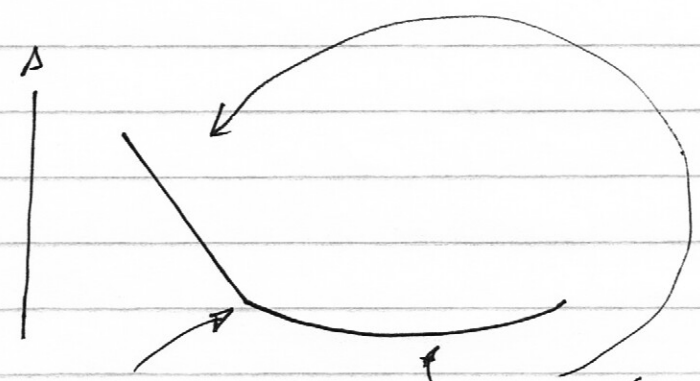
So that I get the same solu.

To proceed, we need a richer notion of gradient, to deal w/ the absolute value.

- the subgradient

Example

convex fn
of one var



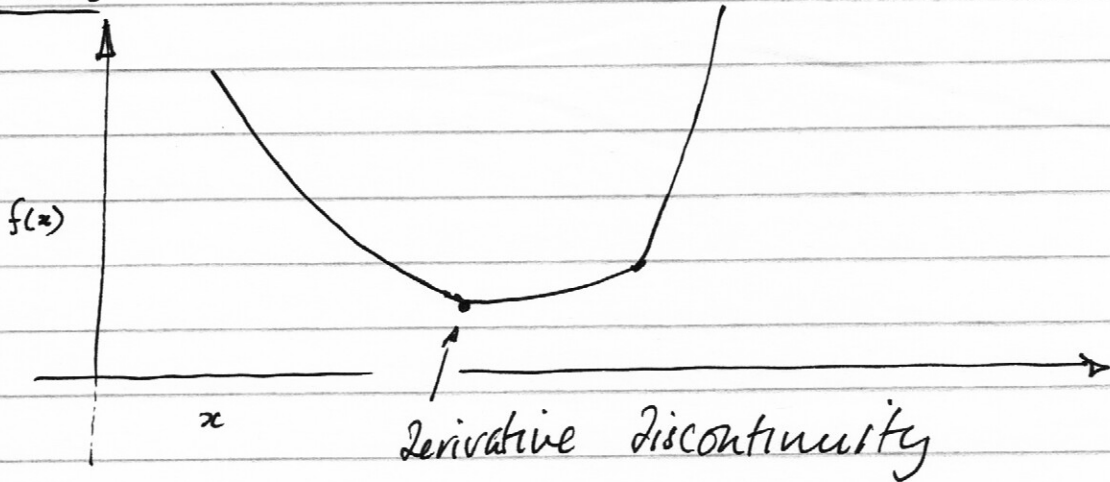
Over here, there is a cone of tangents.

→ Derivative disc.

over here, I can

→ construct unique tangent lines

Subgradients



consider the graph of a convex function
 $(x, f(x))$.

- at a differentiable point, there
 is a tangent plane.

Equation is easy.

$$\text{Surface } g(x) = x_n - f(x_{1:n}) = 0$$

$$\text{Normal} = \nabla g = (-\nabla f, 1)$$

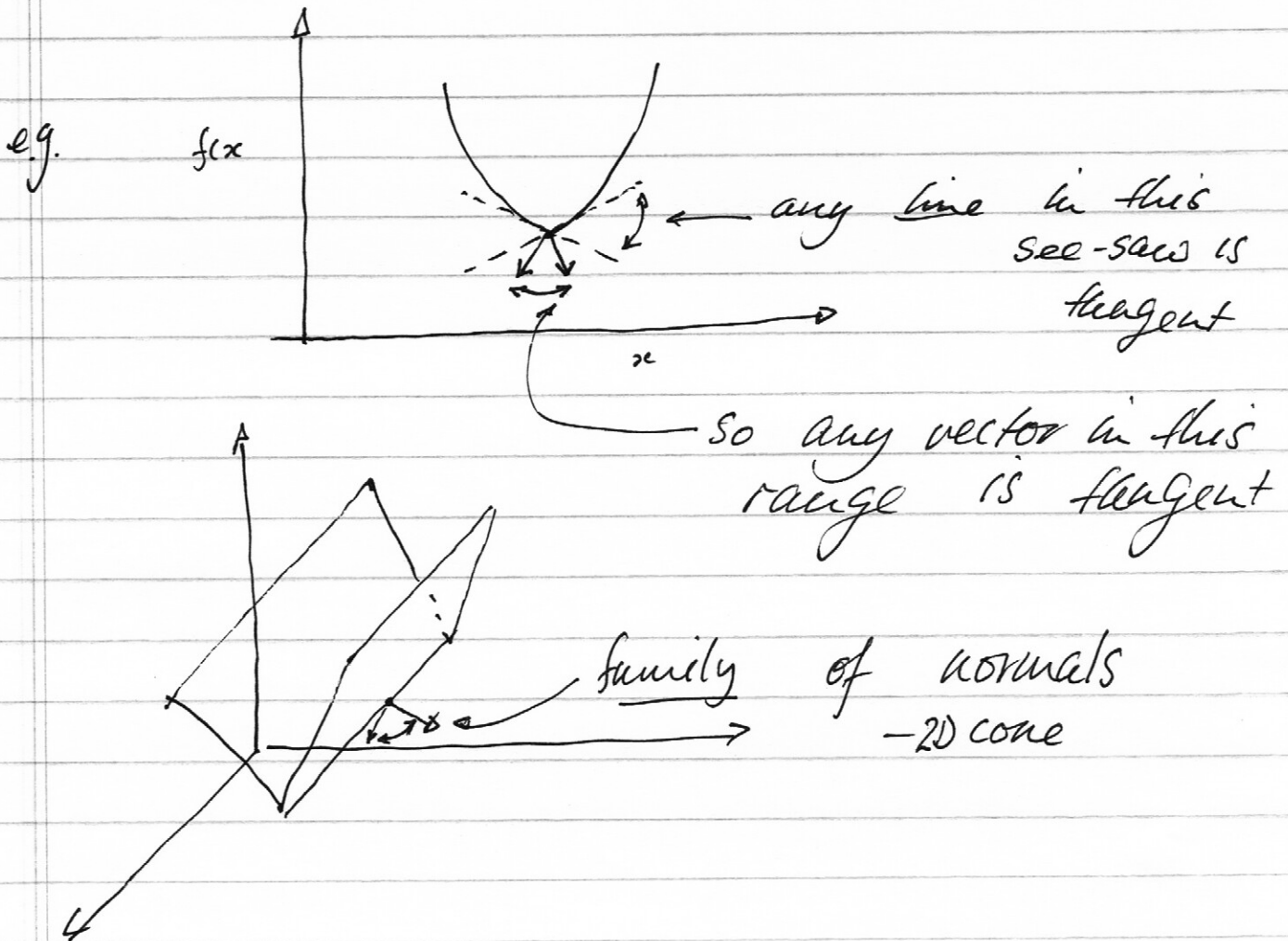
(this isn't unit - doesn't matter)

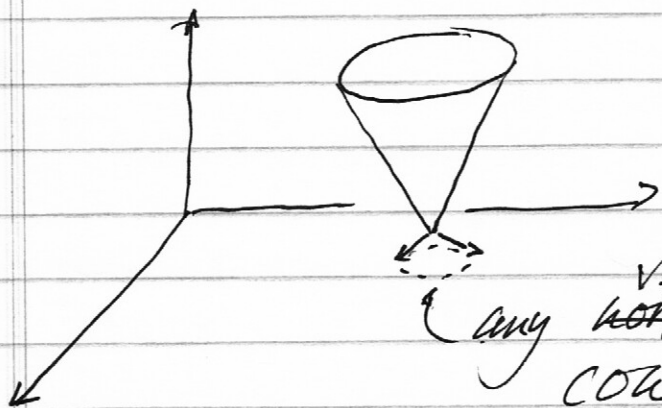
so the tangent plane at
 $(\underline{u}, f(\underline{u}))$

$$\text{is } -\nabla f \Big|_{\underline{u}} \cdot \underline{x}_{1:n} + x_n - \left(-\nabla f \Big|_{\underline{u}} \cdot \underline{u} + f(\underline{u}) \right) = 0$$

Now think about a non-differentiable point.

- there is a cone of normals





any ^{vector} normal in this filled cone is normal.

All this means that at such points, our function has a family of derivatives — known as subgradients

Because a tangent plane (= normal) yields a derivative

TP

Normal

SubGrad

$$-p \cdot x_{1:n} + x_n - \alpha = 0$$

$$(-p, 1)$$

$$p.$$

This yields an easy construction for subgradient of abs

$$\begin{array}{l} \partial |x| = \begin{cases} -1 & x < 0 \\ 1 & x > 0 \\ [-1, 1] & x = 0 \end{cases} \\ \uparrow \\ \text{subgradient} \end{array}$$

↑ closed interval

Because my derivation of ~~Lagrangian~~ derivative condition was geometric, it works for subgradients

So for $f(x)$ to be a min, ~~it~~ must have

$$0 \in \partial f$$

↑ because subgradients produce intervals, or worse

So, for our problem

0

$$0 \in \partial [\|Ax - y\|^2 + \lambda |x|_1]$$

$$= 2A^T(Ax - y) + \lambda \partial |x|_1$$

$$= 2A^T(Ax - y) + \lambda s$$

where $s_i = \begin{cases} 1 & x_i > 0 \\ -1 & x_i < 0 \\ [-1, 1] & x_i = 0 \end{cases}$

s is the sign vector.

Notice if you know the sign vector for a soln, then the soln is easy to get.

(16)

write $J_e(s) = \{i \mid s_i \notin [-1; 1]\}$

$$J_c(s) = \{\text{All indices}\} - J(s)$$

assume we know s ;

then $x_{J_c(s)} = 0$

← wild index notation;
sorry!

so we have

$$A_J^T (A_{JJ} x_J - y_J) + \lambda s_J = 0$$

which is a straightforward linear system (we assume that $A_J^T A_J$ has full rank for any J we deal with - fairly reasonable)

Now, assume

$$\lambda_1 < \lambda_2, \quad S(\lambda_1) = S(\lambda_2) \\ = \sigma$$

$$1) \quad \frac{x(\lambda_1)}{J_c(\sigma)} = \frac{x(\lambda_2)}{J_c(\sigma)} = 0 \quad (\text{obvious})$$

$$2) \quad S(t\lambda_1 + (1-t)\lambda_2) = \sigma \quad (\text{obvious})$$

$$3) \quad A^T \left(A \left[\frac{t x(\lambda_1)}{J_c(\sigma)} + (1-t) \frac{x(\lambda_2)}{J_c(\sigma)} \right] - y \right)_{J_c(\sigma)}$$

$$+ (t\lambda_1 + (1-t)\lambda_2) \sigma = 0$$

(easy)

BUT this means

$x(\lambda)$ is piecewise linear

So it is natural to try and construct
this path. — but there is some
bad news.

write d for dimension of x .

then, clearly, # of verts in
path is $\leq 3^d$

actually, the upper bound is
 $\frac{(3^d + 1)}{2}$, and it can be attained

(see papers).

Algorithms

1) Mathing pursuit

$$r_0 = y; \quad \text{cho } x_0 = 0$$

- choose col j of A that has largest value of $r_i^T A_j$

$$- x_{i+1} = x_i + \frac{(r_i^T A_j)}{(A_j^T A_j)} \cdot e_j$$

$$- r_{i+1} = r_i - \frac{(r_i^T A_j)}{(A_j^T A_j)} \cdot A_j$$

Notice

- residual always gets smaller
- $|x|_1$ gets bigger

Orthogonal matching pursuit

- like matching pursuit, BUT
 - readjust all ~~coeff~~ non-zero coeffs each time you insert a column to get best fit in that space.
 - Better estimate in k steps, but each step takes more work.

Homotopy algorithms

- a whole class of algorithm, constructing approx or exact $x(\lambda)$.

Notice

1) we know $x(0)$

(because $A^T(Ax - y) = 0$, linear alg)

2) for sufficiently large $\lambda = \lambda_{max}$

$$x(\lambda_{max}) = 0$$

- and we can compute λ_{max}

$$0 \in A^T(Ax - y) + \lambda_{max} \begin{bmatrix} -1; 1 \\ \vdots \\ -1; 1 \end{bmatrix}$$

and $x = 0$

so we must have

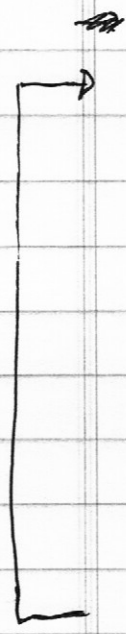
$$0 \in -A^T y + \lambda_{max} \begin{bmatrix} -1 & 1 \\ & \\ & \\ -1 & 1 \end{bmatrix}$$

$$\text{so } \lambda_{max} = \|A^T y\|_{\infty}$$

largest abs. value.

So we can do

$$\lambda_0 = 0 ; \quad \neq \quad A^T A x_0 = A^T y$$



• Tangent to path:

$$\left(\frac{A^T A}{J} \right) \frac{dx_i}{d\lambda} + \frac{\sigma(\lambda)}{J} = 0$$

• search along tangent for i that gives first sign change

• update x

• stop when $\lambda = \lambda_{max}$

Advantage

- complete path
- Disadvantages:

- 1) what if 2 knots coincide (rare)
- 2) too many knots (approximations are available)

Another view of the PL path.

• Problem

$$\min. \|Ax - y\|^2 + \lambda \sum_i t_i$$

$$\text{st } -t_i \leq x_i \leq t_i$$

- Notice this is a family of qp's

- linear constraints, but the polytope

Doesn't change when λ changes

- By inspection, at soln $x_i = \begin{cases} t_i \\ -t_i \end{cases}$

- By inspection, these are 1-faces

or 0-faces

\Rightarrow soln is always on 1-face
or 0-face

But look at SVM (linear)

$$\lambda \frac{w^T w}{2} + \frac{1}{N} \sum_i \xi_i$$

st

$$y_i (w^T x_i + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

1) as λ changes, polytope does not change:

2) at soln, we always have

$$\xi_i = \max(0, 1 - y_i (w^T x_i + b))$$

so a soln is always on a

0-face or 1-face

3) \Rightarrow homotopy path is Ph. !

(and fairly straightforward to construct)