

Matrix differential calculus

10-725 Optimization
Geoff Gordon
Ryan Tibshirani

Review

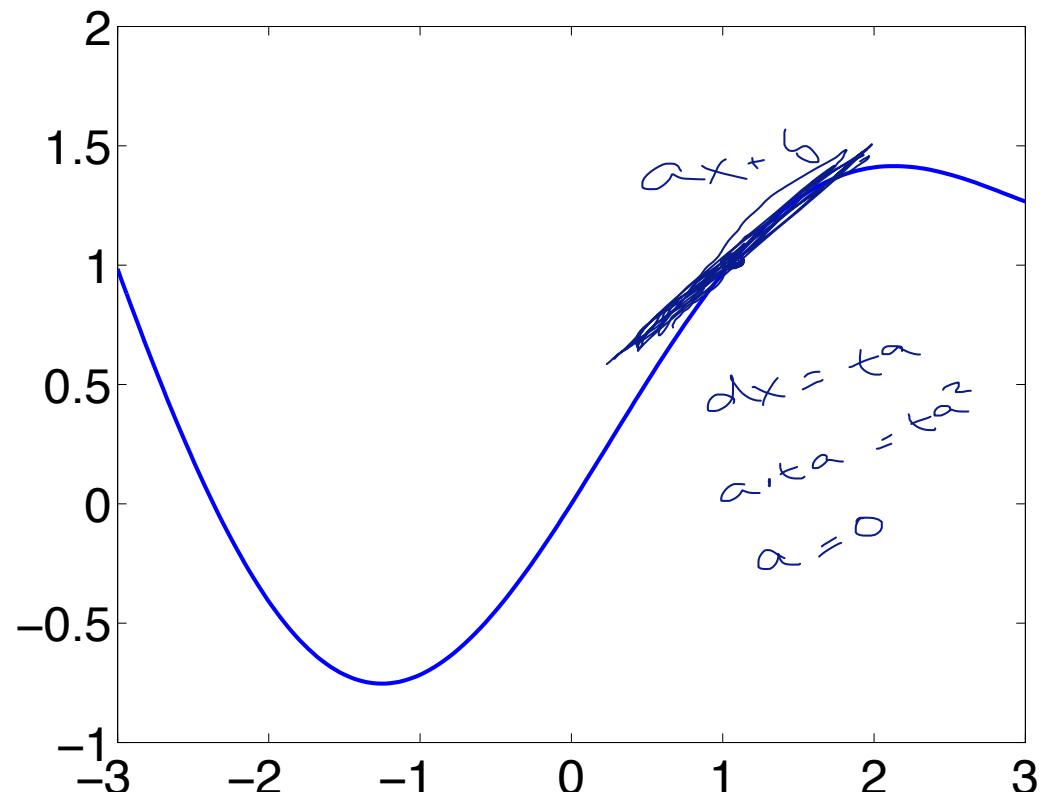
- Matrix differentials: sol'n to matrix calculus pain
 - ▶ compact way of writing Taylor expansions, or ...
 - ▶ definition:
 - ▶ $df = a(x; dx) [+ r(dx)]$
 - ▶ $a(x; .)$ linear in 2nd arg
 - ▶ $r(dx)/\|dx\| \rightarrow 0$ as $dx \rightarrow 0$
- $d(\dots)$ is linear: passes thru +, scalar *
- Generalizes Jacobian, Hessian, gradient, velocity

Review

- Chain rule
- Product rule
- Bilinear functions: cross product, Kronecker, Frobenius, Hadamard, Khatri-Rao, ...
- Identities
 - ▶ rules for working with \circ , $\text{tr}()$
 - ▶ trace rotation $\text{tr}(AB) = \text{tr}(BA)$
- Identification theorems

Finding a maximum or minimum, or saddle point

ID for $df(x)$	scalar x	vector \mathbf{x}	matrix X
scalar f	$df = a \, dx$	$df = \mathbf{a}^T d\mathbf{x}$	$df = \text{tr}(A^T dX)$
vector \mathbf{f}	$Df = A \, dx$	$Df = \nabla f(\mathbf{x})^T d\mathbf{x}$	$Df = \nabla f(X)^T dX$

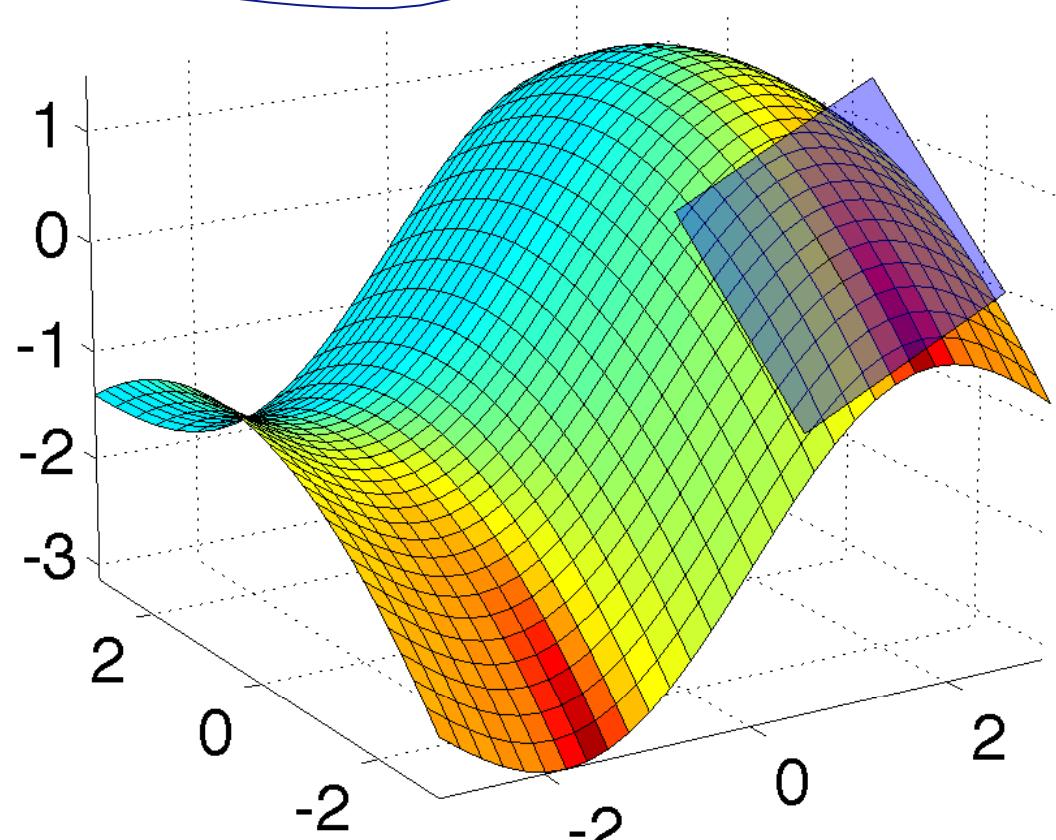


Finding a maximum

or minimum, or saddle point

ID for $df(x)$	scalar x	vector \mathbf{x}	matrix \mathbf{X}
scalar f	$df = a \, dx$	$df = \mathbf{a}^T d\mathbf{x}$	$df = \text{tr}(\mathbf{A}^T d\mathbf{X})$

$$dx = t \vec{a}$$
$$df \approx t \vec{a}^T \vec{a} > 0$$
$$f \text{ at } t=0$$



And so forth...

- Can't draw it for X a matrix, tensor, ...
- But same principle holds: set coefficient of dX to 0 to find min, max, or saddle point:
 - ▶ if $df = c(A; dX) [+ r(dX)]$ then

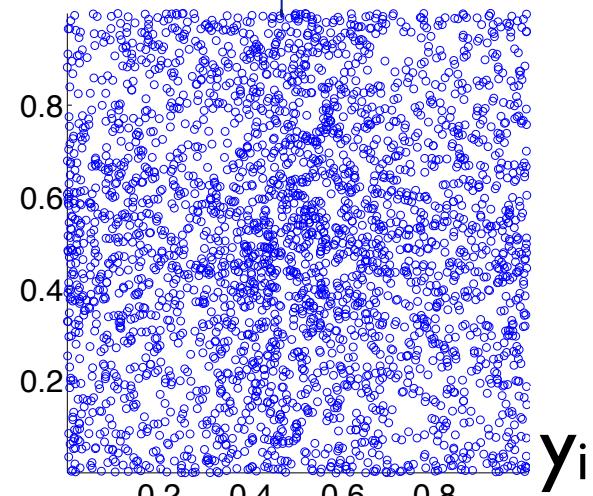
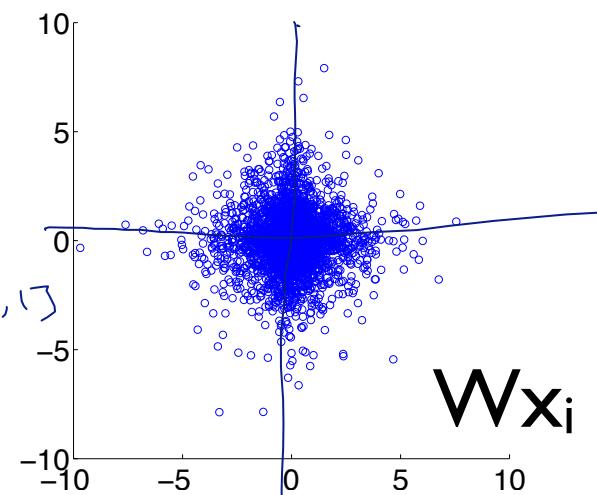
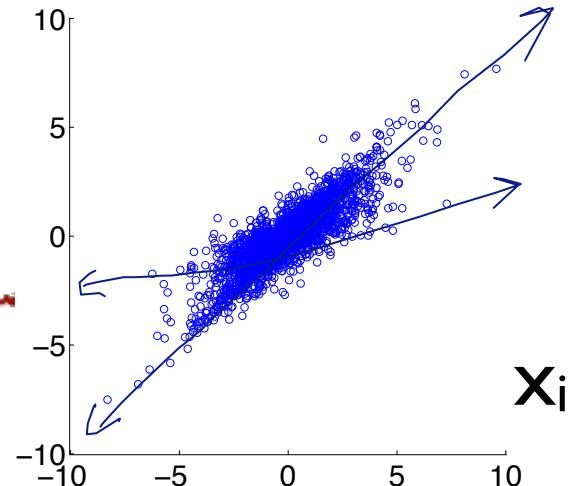
$$dX = t \hat{A} \quad dF = c(A; t \hat{A}) = t \underline{c(A; A)}$$

- ▶ so: max/min/sp iff $c(A; \hat{A}) = 0$
- ▶ for $c(\cdot, \cdot)$ any “product”, $\Leftrightarrow \hat{A} = 0$

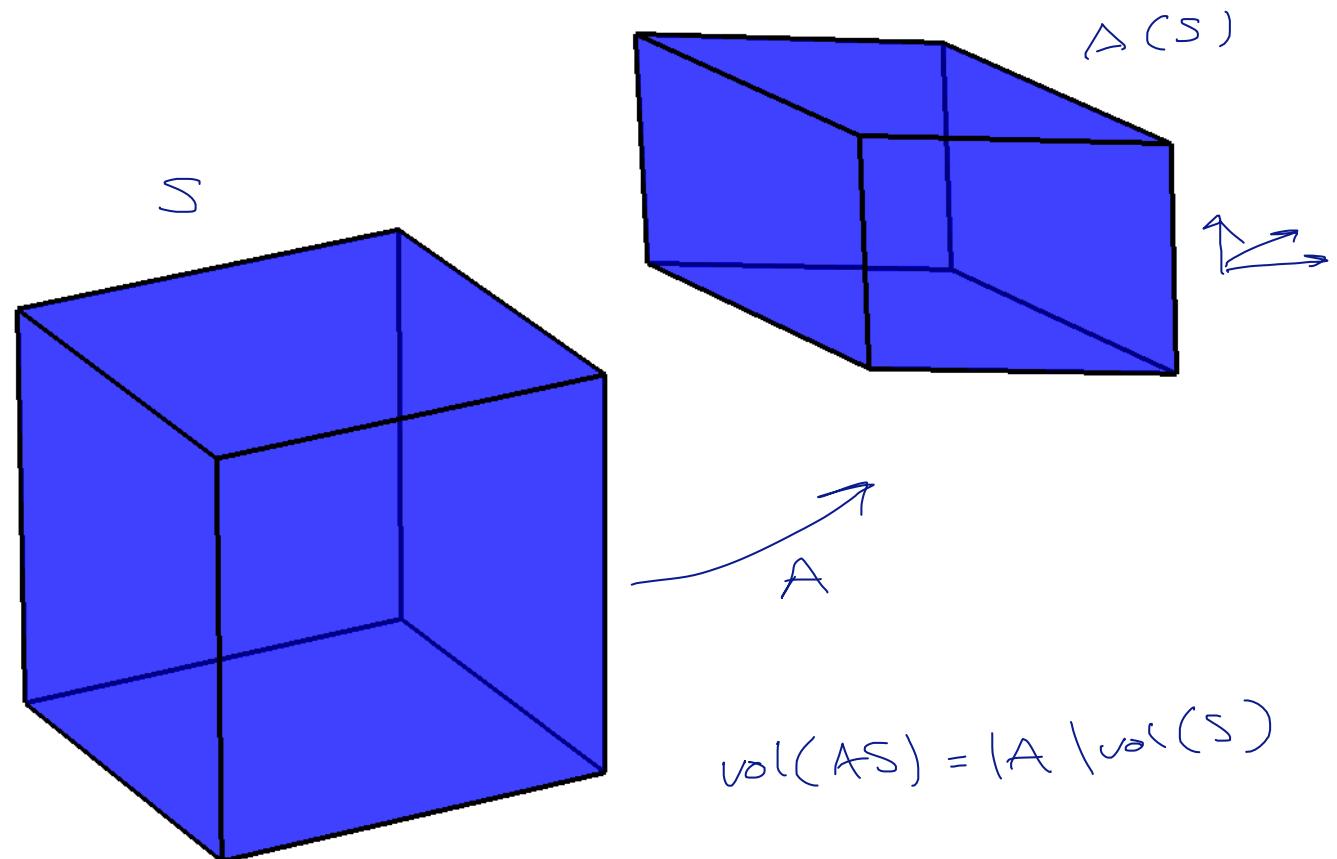
Ex: Infomax ICA

- Training examples $x_i \in \mathbb{R}^d, i = 1:n$
- Transformation $y_i = g(Wx_i)$
 - ▶ $W \in \mathbb{R}^{d \times d}$ parameter
 - ▶ $g(z) = \text{scalar fn, componentwise } g: \mathbb{R} \rightarrow [0, 1]$
- Want: independent components

$\max F \rightarrow$



Volume rule



Ex: Infomax ICA

- $y_i = g(Wx_i)$

► $dy_i = J(x_i, w)dx_i = J_i dx_i = -\underline{\mathbb{E} \langle \ln P(y_i) \rangle}$

$$\max H(P(y))$$

$$= - \int P(y) \ln P(y) dy$$

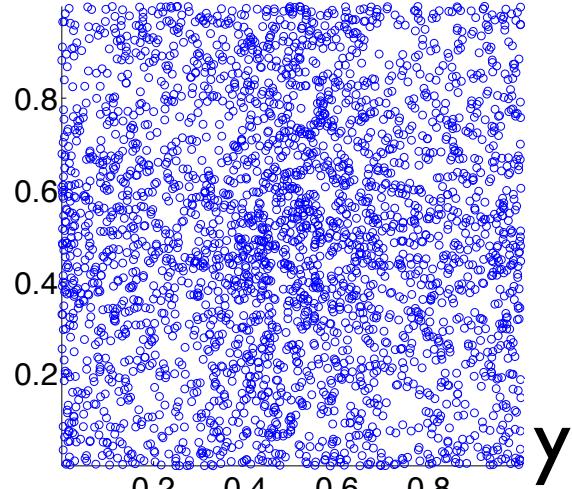
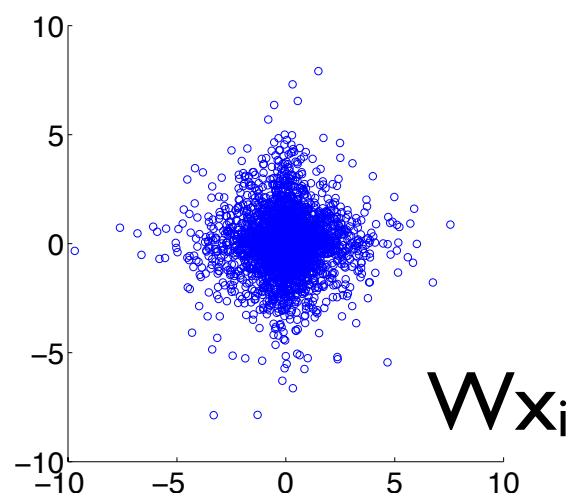
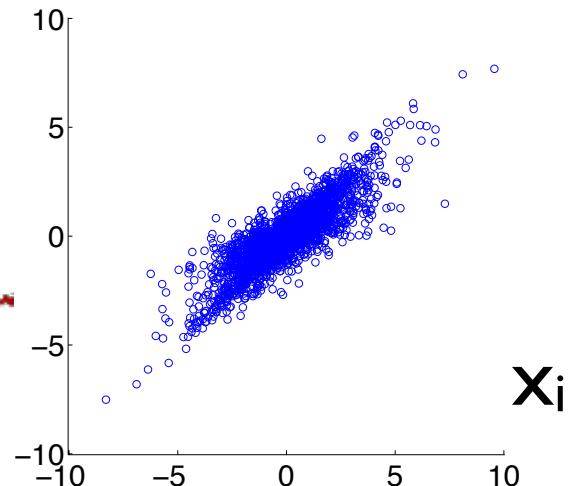
$$= -\mathbb{E} \langle \ln P(y) \rangle$$

\downarrow est. of $-\int P(y_i) \ln P(y_i)$

- Method: $\max_w \sum_i -\ln(P(y_i))$

► where $P(y_i) = P(x_i) / |\det J(x_i, w)|$

$$\max_w \sum_i (\ln |\det J(x_i, w)| - \ln P(x_i))$$



Gradient

$$J_i = \text{diag}(u_i) \omega$$

$$\bullet L = \sum_i \ln |\det J_i| \quad y_i = g(Wx_i) \quad dy_i = J_i d\omega$$

$$\begin{aligned} dJ_i &= (\underbrace{d \text{diag}(u_i)}_{\text{diag}(u_i) d\omega} \omega \\ &\quad + \text{diag}(u_i) d\omega \\ &= \text{diag}(v_i d\omega x_i) \omega \\ &\quad + \text{diag}(u_i) d\omega \end{aligned}$$

$$\begin{aligned} dy_i &= g'(Wx_i) \circ d(Wx_i) \quad u_i = g'(Wx_i) \\ &\leftarrow u_i \circ (Wdx_i) \quad du_i = \underbrace{g''(Wx_i)}_{v_i} \circ (dWx_i) \\ &= \text{diag}(u_i) \omega dx_i \quad d \text{diag}(u_i) = \text{diag}(v_i = dWx_i) \end{aligned}$$

Gradient

$$J_i = \text{diag}(u_i) W \quad dJ_i = \text{diag}(u_i) dW + \text{diag}(v_i) \text{diag}(dW x_i) W$$

$$\begin{aligned} dL &= \sum_i d \ln |\det J_i| \\ &= \sum_i \text{tr} (\tilde{J}_i^{-1} d\tilde{J}_i) \\ &= \sum_i \text{tr} (\omega^{-1} d\omega + \underbrace{\omega^{-1} \text{diag}(u_i)^{-1} \text{diag}(u_i) \text{diag}(d\omega x_i)}_{\text{diag}(\alpha_i)}) \\ &= \sum_i \text{tr} (\omega^{-1} d\omega) + \underbrace{\text{tr} \text{diag}(\alpha_i) \text{diag}(d\omega x_i)}_{x_i^\top d\omega x_i} \\ &\sim \underbrace{\sum_i \text{tr} (\omega^{-1} d\omega)}_{n \omega^{-1} d\omega} + \text{tr} ((\sum_i x_i \alpha_i^\top) d\omega) \\ \frac{dL}{d\omega} &= n \omega^{-T} + \sum_i \alpha_i x_i^\top = n (\omega^{-T} + C) \end{aligned}$$

Natural gradient

- $L(W): \mathbb{R}^{d \times d} \rightarrow \mathbb{R}$ $dL = \text{tr}(G^T dW)$
- step $S = \arg \max_S M(S) = \text{tr}(G^T S) - \|S W^{-1}\|_F^2 / 2$
 - scalar case: $M = gs - s^2 / 2w^2$
- $M = \text{tr}(G^T S) - \frac{1}{2} \text{tr}(S \omega^{-1} \omega^{-T} S^T)$
- $dM = \text{tr}(G^T dS) - \text{tr}(\underbrace{\omega S \omega^{-1} \omega^{-T} S^T}_{\circlearrowright})$

$$G = S \omega^{-1} \omega^{-T}$$

$$G \omega^T \omega = S$$

ICA natural gradient

- $[W^T + C] W^T W = \text{I} + C w^T w$

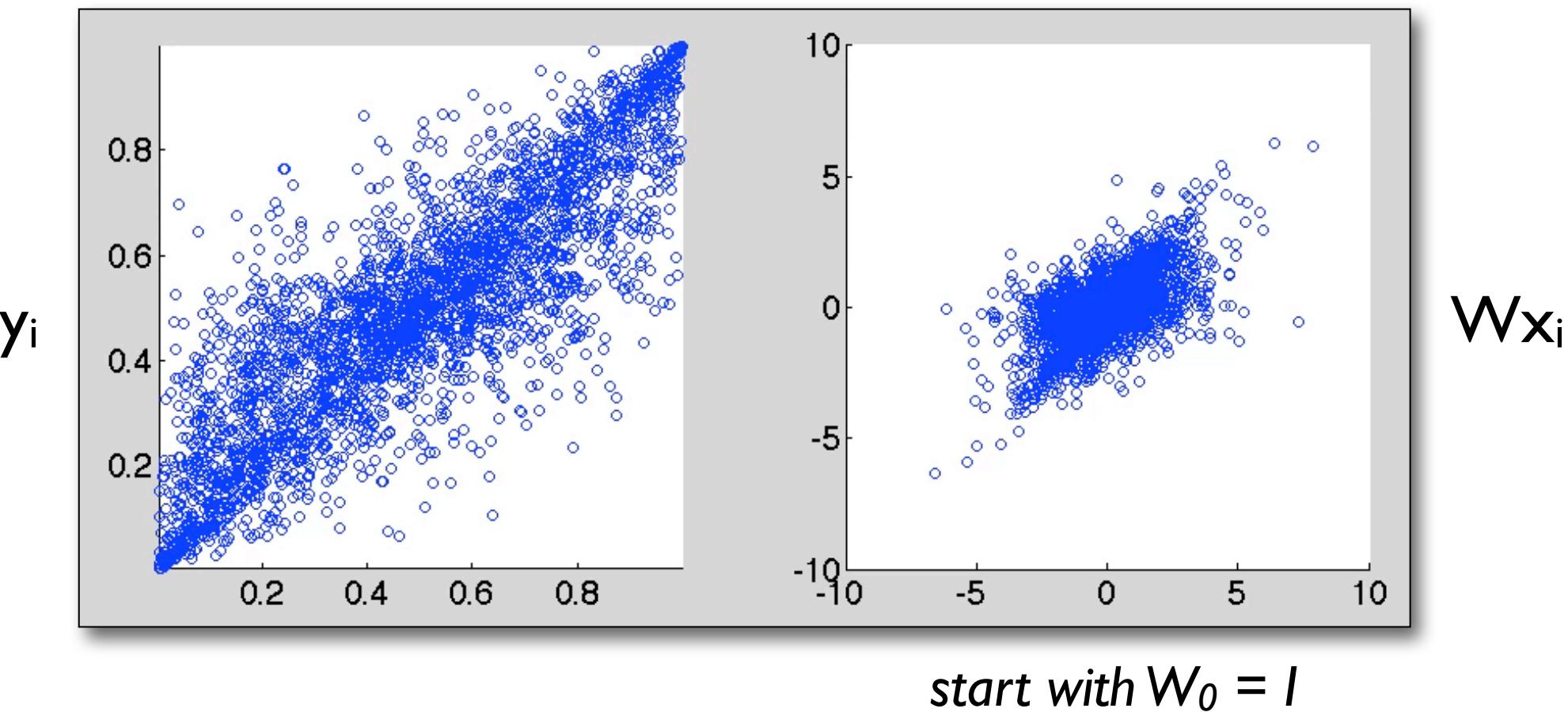
y_i

Wx_i

start with $W_0 = I$

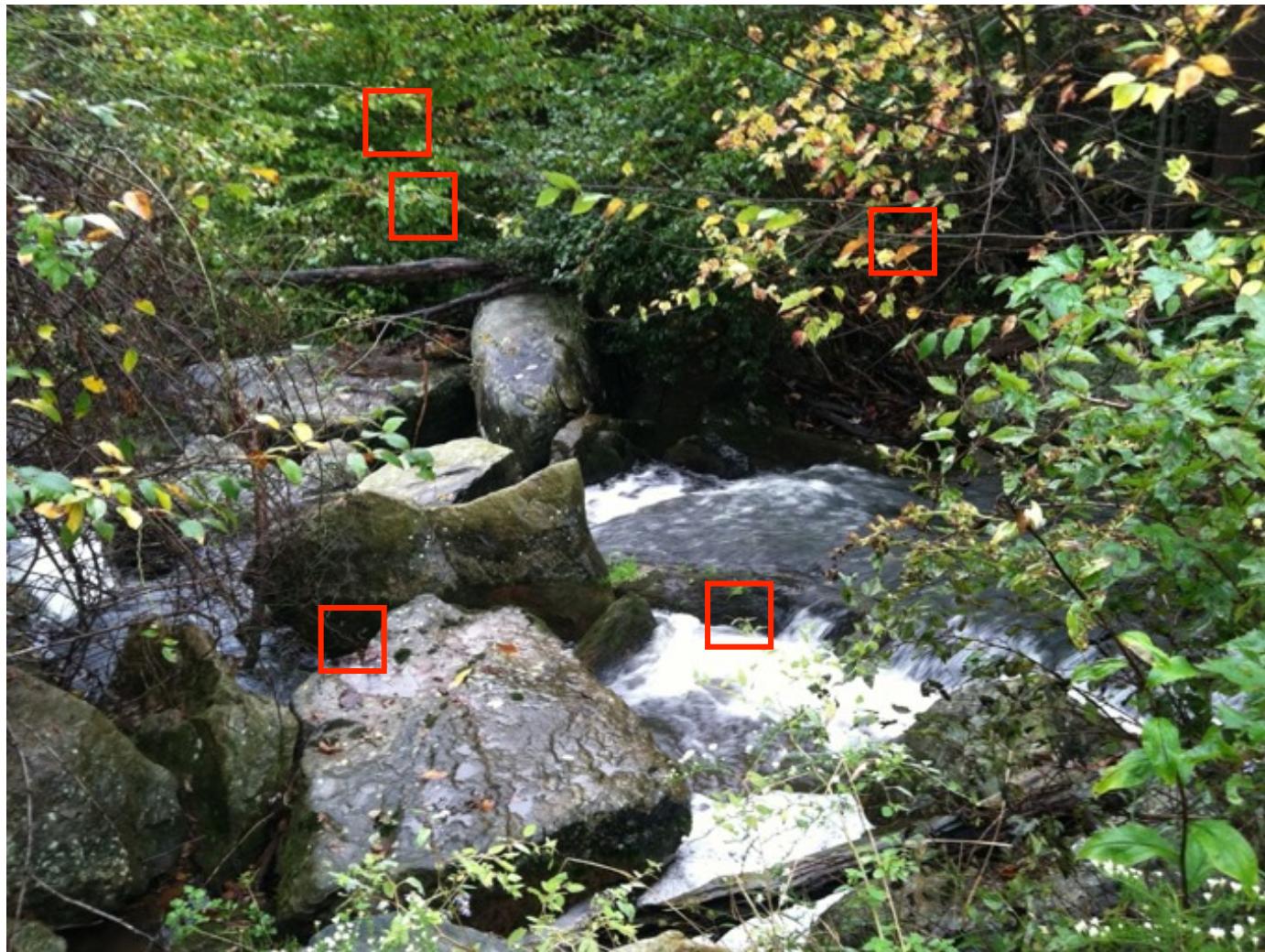
ICA natural gradient

- $[W^{-T} + C] W^T W =$

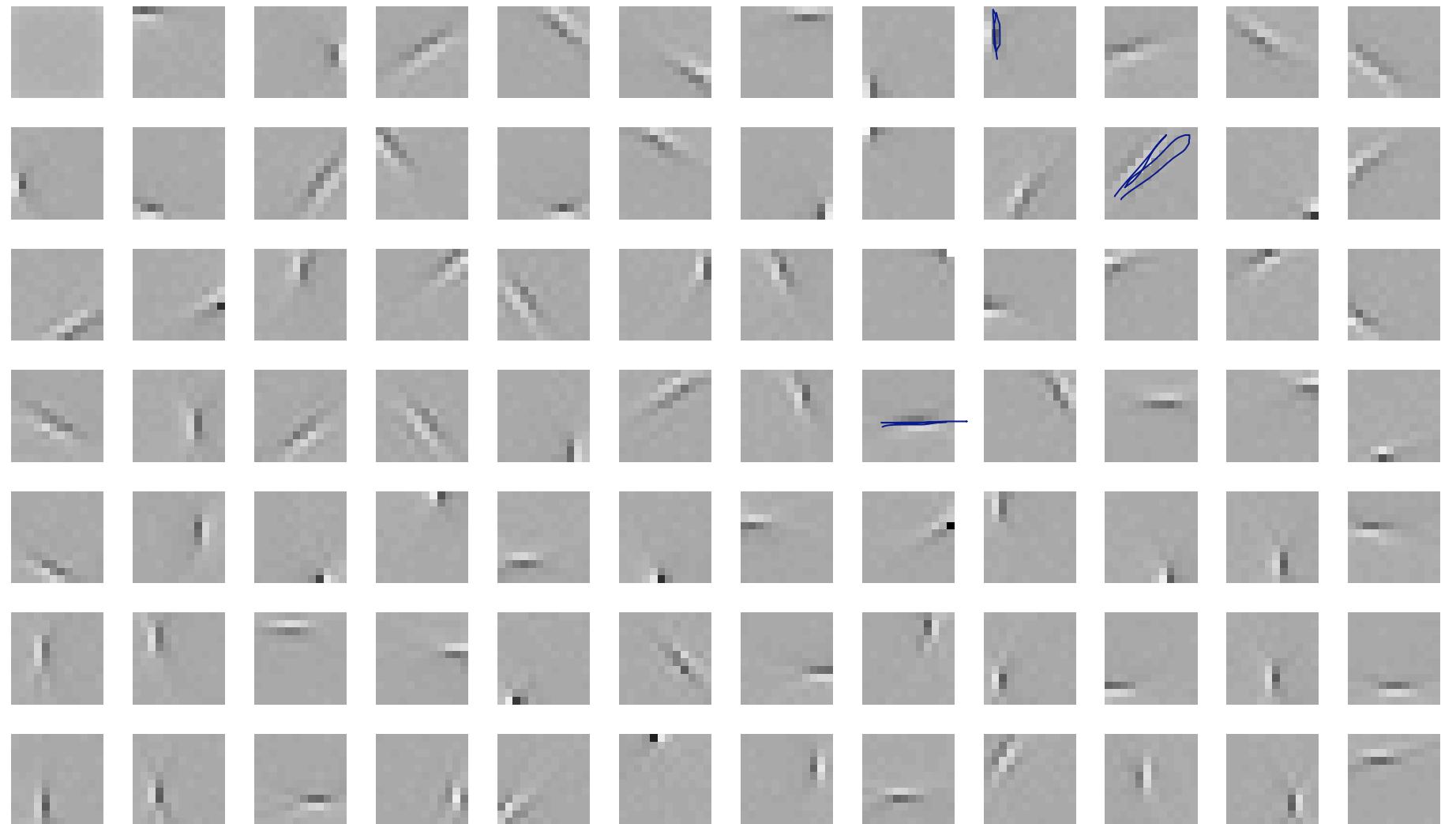


start with $W_0 = I$

ICA on natural image patches



ICA on natural image patches



More info

- Minka's cheat sheet:
 - ▶ <http://research.microsoft.com/en-us/um/people/minka/papers/matrix/>
- Magnus & Neudecker. *Matrix Differential Calculus*. Wiley, 1999. 2nd ed.
 - ▶ <http://www.amazon.com/Differential-Calculus-Applications-Statistics-Econometrics/dp/047198633X>
- Bell & Sejnowski. An information-maximization approach to blind separation and blind deconvolution. *Neural Computation*, v7, 1995.

Newton's method



10-725 Optimization
Geoff Gordon
Ryan Tibshirani

Nonlinear equations

- $x \in \mathbb{R}^d \quad f: \mathbb{R}^d \rightarrow \mathbb{R}^d$, diff'ble

► solve: $f(x) = 0$

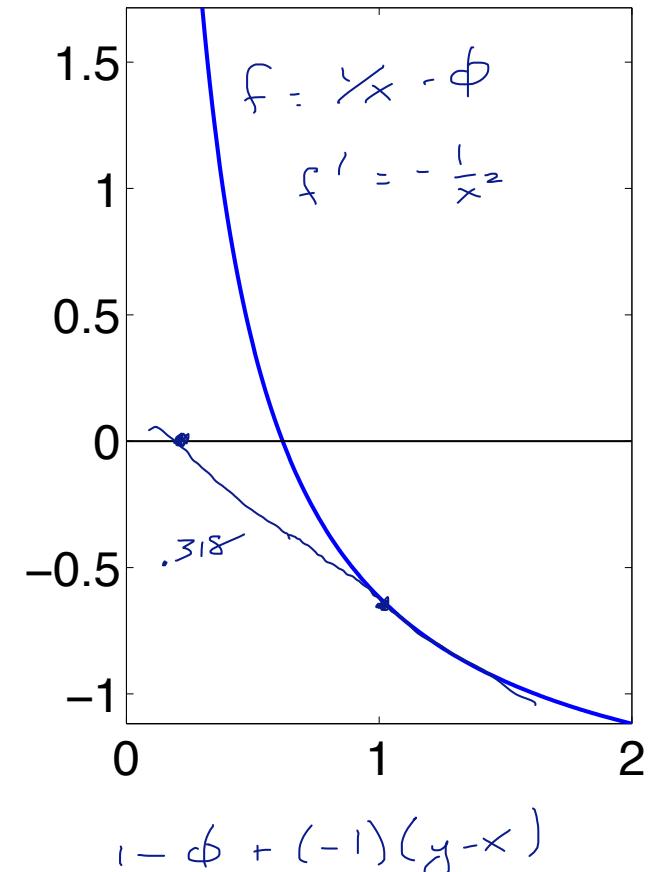
- Taylor: $f(y) \approx f(x) + J(x)(y - x) = \hat{f}(y)$

► J : Jacobian

- Newton: $\hat{f}(y) = 0$

$$f(x) + J(x)(y - x) = 0$$

$$y - x = -J(x)^{-1} f(x)$$



$$\hat{f} = \frac{f(x) - \frac{1}{x^2} dx}{1/x - \phi} = 0$$

$$\frac{1 - \phi - \frac{1}{x}}{x - x^2 \phi - dx} = 0$$

Error analysis

$$\epsilon = x\phi - 1 \quad x^+ = x + x(1-x\phi)$$

$$\begin{aligned}\epsilon^+ &= x^+\phi - 1 \\ &= [x + \underbrace{x(1-x\phi)}_{-\epsilon}] \phi - 1\end{aligned}$$

$$= x\phi - x\epsilon\phi - 1$$

$$\begin{aligned}&= +\epsilon - x\epsilon\phi \\ &= \epsilon(1 - x\phi) = -\epsilon^2\end{aligned}$$

$$dx = x * (1 - x * \phi)$$

0:	0.7500000000000000
1:	0.5898558813281841
2:	0.6167492604787597
3:	0.6180313181415453
4:	0.6180339887383547
5:	0.6180339887498948
6:	0.6180339887498949
7:	0.6180339887498948
8:	0.6180339887498949
<hr/>	
*	0.6180339887498948

Bad initialization

1.300000000000000
-0.1344774409873226
-0.2982157033270080
-0.7403273854022190
-2.3674743431148597
-13.8039236412225819
-335.9214859516196157
-183256.0483360671496484
-54338444778.1145248413085938

Minimization

- $x \in \mathbb{R}^d \quad f: \mathbb{R}^d \rightarrow \mathbb{R}$, twice diff'ble

► find: $\min_x f(x) \quad o = f'(x) \doteq g(x)$

$$g: \mathbb{R}^d \rightarrow \mathbb{R}^d$$

- Newton:

$$\begin{aligned} dx &= - (f'(x))^{-1} g(x) \\ &= - (f''(x))^{-1} f'(x) \\ &= - H^{-1} g \end{aligned}$$

Descent

- Newton step: $d = -(f''(x))^{-1} f'(x)$
- Gradient step: $-g = -f'(x)$
- Taylor: $df = g^T dx \quad [+ r(dx)]$
- Let $t > 0$, set $dx = t d$
 - ▶ $df = \underbrace{-t f'(x) (f''(x))^{-1} f'(x)}_{> 0 \text{ iff } f''(x) \succ 0} + r(td)$
- So:

Steepest descent

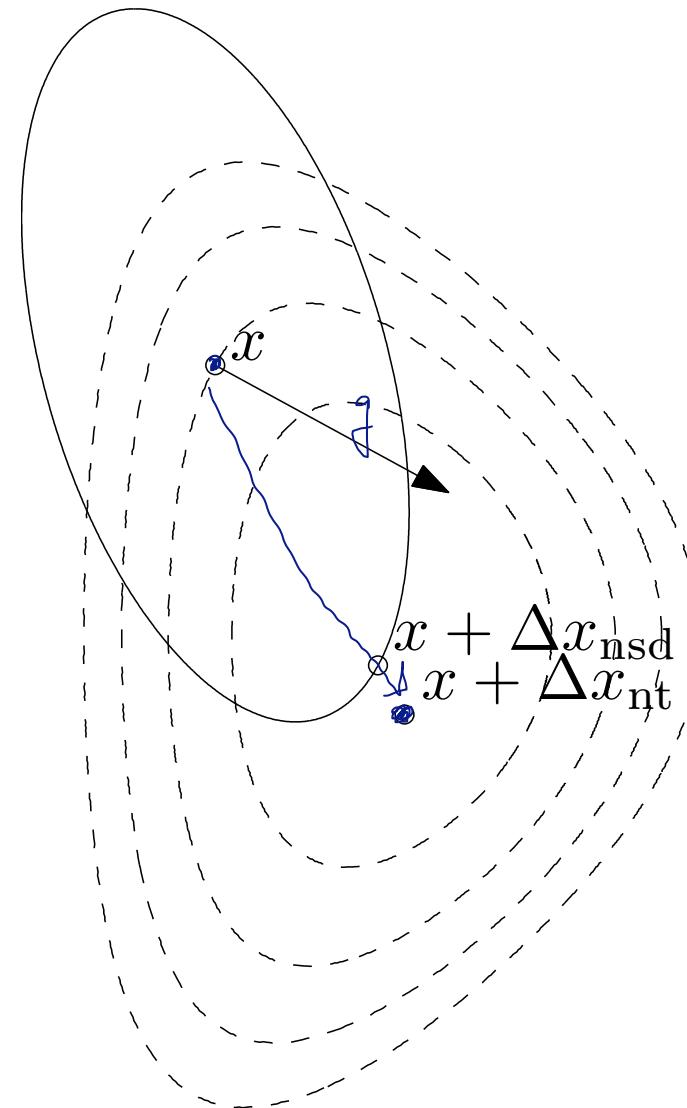
$$g = f'(x)$$

$$H = f''(x)$$

$$\|d\|_H = \sqrt{d^T H d}$$

$$\min_d g^T d + \underbrace{\|d\|_H^2}_{d^T H d}$$

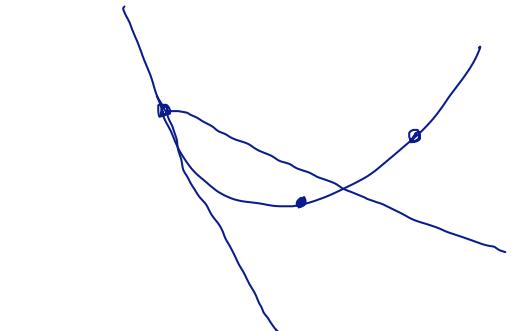
$$d = -H^{-1}g$$



"damped"

Newton w/ line search

- Pick x_1
- For $k = 1, 2, \dots$
 - ▶ $g_k = f'(x_k); H_k = f''(x_k)$
 - ▶ $d_k = -H_k \setminus g_k$
 - ▶ $t_k = 1$
 - ▶ while $f(x_k + t_k d_k) > f(x_k) + t \underline{g_k^T d_k / 2}$
 - ▶ $t_k = \beta t_k$ $\beta < 1$
 - ▶ $x_{k+1} = x_k + t_k d_k$ step



gradient & Hessian

Newton direction

backtracking line search

Properties of damped Newton

- Affine invariant: suppose $g(x) = f(Ax+b)$
 - ▶ x_1, x_2, \dots from Newton on $g()$
 - ▶ y_1, y_2, \dots from Newton on $f()$
 - ▶ If $y_1 = Ax_1 + b$, then: $y_i = Ax_i + b \quad \forall i$
- Convergent:
 - ▶ if f bounded below, $f(x_k)$ converges
 - ▶ if f strictly convex, bounded level sets, x_k converges
 - ▶ typically quadratic rate in neighborhood of x^*