

# Optimization for well-behaved problems

For statistical learning problems, “well-behaved” means:

- signal to noise ratio is decently high
- correlations between predictor variables are under control
- number of predictors  $p$  can be larger than number of observations  $n$ , but not absurdly so

For well-behaved learning problems, people have observed that gradient or generalized gradient descent can converge extremely quickly (much more so than predicted by  $O(1/k)$  rate)

Largely unexplained by theory, topic of current research. E.g., very recent work<sup>4</sup> shows that for some well-behaved problems, w.h.p.:

$$\|x^{(k)} - x^\star\|^2 \leq c^k \|x^{(0)} - x^\star\|^2 + o(\|x^\star - x^{\text{true}}\|^2)$$

---

<sup>4</sup>Agarwal et al. (2012), *Fast global convergence of gradient methods for high-dimensional statistical recovery*

# *Administrivia*



- HW2 out as of this past Tuesday—due 10/9
- Scribing
  - ▶ Scribes 1–6 ready soon; handling errata
  - ▶ missing days: 11/6, 12/4, and 12/6
- Projects:
  - ▶ you should expect to be contacted by TA mentor in next weeks
  - ▶ project milestone: 10/30

# Matrix differential calculus



*10-725 Optimization*  
*Geoff Gordon*  
*Ryan Tibshirani*

# *Matrix calculus pain*



- Take derivatives of fns involving matrices:
  - ▶ write as huge multiple summations w/ lots of terms
  - ▶ take derivative as usual, introducing more terms
    - ▶ case statements for  $i = j$  vs.  $i \neq j$
  - ▶ try to recognize that output is equivalent to some human-readable form
  - ▶ hope for no indexing errors...
- Is there a better way?

# Differentials

- Assume  $f$  sufficiently “nice”
- Taylor:  $f(y) = f(x) + f'(x)(y-x) + r(y-x)$   
→ with  $r(y-x) / |y-x| \rightarrow 0$  as  $y \rightarrow x$   $x$  fixed
- Notation:  $df = f(y) - f(x)$      $dx = y - x$     “differentials”  
 $df = f'(x)dx [ + r(dx) ]$

# Definition

- Write

- ▶  $dx = y - x$

- ▶  $df = f(y) - f(x)$

- Suppose

- ▶  $df = a(x; dx) + r(dx)$

- ▶ with  $a$  linear in  $dx$

- ▶ and  $r(dx)/|dx| \rightarrow 0$  as  $dx \rightarrow 0$

- Then:  $a(x; dx)$  is differential of  $f$

$$\begin{aligned} a(x; \lambda dx) &= \lambda a(x; dx) \\ a(x; dx_1 + dx_2) &= a(x; dx_1) + a(x; dx_2) \end{aligned}$$

# Matrix differentials

- For matrix  $X$  or matrix-valued function  $F(X)$ :

▶  $dX = [dx_{ij}]_{ij}$

▶  $dF = a(x; dX) + r(dX)$

▶ where  $a$  linear in  $dX$

▶ and  $r(dX)/\|dX\| \rightarrow 0$  as  $dX \rightarrow 0$

- Examples:  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$   $df = J(x) dX$

$$f: \mathbb{R}^n \rightarrow \mathbb{R} \quad df = g(x)^T dx \quad dg = H(x) dx$$

$\uparrow \mathbb{R}^n$ 
 $\uparrow \mathbb{R}^{n \times n}$

# Working with differentials

- Linearity:

- ▶  $d(f(x) + g(x)) = df(x) + dg(x)$

- ▶  $d(k f(x)) = k df(x)$

- If  $g$  linear,  $dg(f(x)) = g(df(x))$ ; for example

- ▶  $\text{reshape}(A, [m \ n \ k \ \dots])$

- ▶  $\text{vec}(A) = A(:) = \text{reshape}(A, [], 1)$

- ▶  $\text{tr}(A) = \sum_i A_{ii}$

- ▶  $A^T$



# Reshape

```
>> A = reshape(1:24, [2 3 4])
```

```
A(:, :, 1) =
```

1	3	5
2	4	6

```
A(:, :, 2) =
```

7	9	11
8	10	12

```
A(:, :, 3) =
```

13	15	17
14	16	18

```
A(:, :, 4) =
```

19	21	23
20	22	24

```
>> B = reshape(A, [4 3 2])
```

```
B(:, :, 1) =
```

1	5	9
2	6	10
3	7	11
4	8	12

```
B(:, :, 2) =
```

13	17	21
14	18	22
15	19	23
16	20	24

# Working with differentials

- Chain rule:  $L(x) = f(g(x))$

► want:  $dL(x) = f'(g(x))g'(x)dx$

► have:  $df = a(g(x); dg) [ + r(dg) ]$   
 $dg = b(x; dx) [ + s(dx) ]$

$$r(dg)/\|dg\| \rightarrow 0 \text{ as } dg \rightarrow 0$$

$$\begin{aligned} &\rightarrow a(g(x); dg) \rightarrow f'(g(x))dg \\ &\rightarrow g'(x)dx \end{aligned}$$

$$\begin{aligned} dL &= a(g(x); b(x; dx) + s(dx)) + r(dg) \\ &= \underline{a(g(x); b(x; dx))} + [a(g(x); s(dx)) + r(dg)] \end{aligned}$$

$$\begin{aligned} &a(g(x); s(dx))/\|dx\| + r(dg)/\|dx\| \rightarrow 0 \text{ as } dx \rightarrow 0 \\ &a(g(x); \frac{s(dx)}{\|s(dx)\|}) \frac{\|s(dx)\|}{\|dx\|} + r(dg)/\|dg\| \frac{\|dg\|}{\|dx\|} \end{aligned}$$

# Working with differentials

- Product rule:  $L(x) = c(f(x), g(x))$ 
  - ▶ where  $c$  is **bilinear** = linear in each argument (with other argument fixed)
  - ▶ e.g.,  $L(x) = f(x)g(x)$ :  $f, g$  scalars, vectors, or matrices

$$dL = c(df; g(x)) + c(f(x); dg)$$

# Lots of products

- Cross product:  $d(a \times b) = da \times b + a \times db$
- Hadamard product  $A \circ B = A .* B$   $\leftarrow$  commutative
  - ▶  $(A \circ B)_{ij} = A_{ij} B_{ij}$
  - ▶  $d(A \circ B) = dA \circ B + A \circ dB$
- Kronecker product  $d(A \otimes B) = dA \otimes B + A \otimes dB$
- Frobenius product  $A:B = \sum_{ij} A_{ij} B_{ij}$
- Khatri-Rao product:  $d(A * B) = dA * B + A * dB$

# Kronecker product

```
>> A = reshape(1:6, 2, 3)
```

```
A =
```

1	3	5
2	4	6

```
>> B = 2*ones(2)
```

```
B =
```

2	2
2	2

```
>> kron(A, B)
```

```
ans =
```

2	2	6	6	10	10
2	2	6	6	10	10
4	4	8	8	12	12
4	4	8	8	12	12

```
>> kron(B, A)
```

```
ans =
```

2	6	10	2	6	10
4	8	12	4	8	12
2	6	10	2	6	10
4	8	12	4	8	12

# Hadamard product

- a, b vectors

- ▶  $a \circ b = \text{diag}(a) b = \text{diag}(b) a$

- ▶  $\text{diag}(a) \text{diag}(b) = \text{diag}(a \circ b)$

- ▶  $\text{tr}(\text{diag}(a) \text{diag}(b)) = a^T b$

- ▶  $\text{tr}(\text{diag}(b)) = \text{tr}(\text{diag}(1) \overset{v}{\text{diag}(b)}) = 1^T b$

# Some examples

- $L = (Y - XW)^T(Y - XW)$ : differential wrt  $W$

$$\begin{aligned}\triangleright dL &= d[Y^T Y - Y^T X W - W^T X^T Y + W^T X^T X W] \\ &= -Y^T X dW - dW^T X^T Y + dW^T X^T X W + W^T X^T X dW\end{aligned}$$

# Some examples

matrix  $X$

- $L = [x_{ij}]^2 / 2$        $dL = \begin{bmatrix} 0 & 0 & \dots & 0 \\ \vdots & x_{ij} & \vdots & 0 \\ 0 & \dots & 0 & 0 \end{bmatrix} \circ dX$

- $L = [x^T x]_{ij}^2 / 2$

►  $dL = \begin{bmatrix} 0 & 0 & \dots & 0 \\ \vdots & 2x_{ij} & \vdots & 0 \\ 0 & \dots & 0 & 0 \end{bmatrix} \circ d(x^T x)$

$$= M \circ (dX^T X + X^T dX)$$



# Trace

- $\text{tr}(A) = \sum_i A_{ii}$

- ▶  $d \text{tr}(f(x)) = \text{tr}(d f(x))$

scalar

- ▶  $\text{tr}(x) = x$

- ▶  $\text{tr}(X^T) = \text{tr}(X)$

- Frobenius product:

- ▶  $A:B = \sum_{ij} A_{ij} B_{ij} = \text{tr}(A^T B)$

$$\sum_k A_{ki} B_{kj}$$

$$(A^T B)_{ij} = \sum_k A_{ki} B_{kj}$$

$$\sum_i (A^T B)_{ii} = \sum_i \sum_k A_{ki} B_{ki}$$

$$\sum_i \sum_k A_{ki} B_{ki}$$

# Trace rotation

- $\text{tr}(AB) = \text{tr}(BA)$
- $\text{tr}(ABC) = \text{tr}(CAB) = \text{tr}(BCA)$ 
  - ▶  $\text{size}(A)$ :  $m \times n$
  - ▶  $\text{size}(B)$ :  $n \times \cancel{k} \leftarrow k$
  - ▶  $\text{size}(C)$ :  $k \times m$

# More

- Identities: for a matrix  $X$ ,
  - ▶  $d(X^{-1}) = -X^{-1} (dX) X^{-1}$
  - ▶  $d(\det X) = (\det X) \operatorname{tr}(X^{-1} dX)$
  - ▶  $d(\ln |\det X|) = \operatorname{tr}(X^{-1} dX)$
  - ▶ ...

# Example: linear regression

- Training examples:  $(x_i, y_i) \quad i = 1, \dots, n$

- Input feature vectors:  $x_i \in \mathbb{R}^d$

- Target vectors:  $y_i \in \mathbb{R}^k$

- Weight matrix:  $W \in \mathbb{R}^{k \times d}$

- $\min_W L = \sum_i \|y_i - Wx_i\|_2^2$

► as matrix:  $\|Y - WX\|_F^2$

The diagram shows the matrix equation  $Y = WX$  with dimensions indicated. Matrix  $Y$  is  $k \times n$ , matrix  $W$  is  $k \times d$ , and matrix  $X$  is  $d \times n$ . The result  $Y$  is  $k \times n$ .

I accidentally transposed  $WX \rightarrow XW$  here, compared to the previous slide -- so for this slide only, the regression is from rows of  $X$  to rows of  $Y$

# Linear regression

can do it either way...

$$\bullet L = \|Y - \overset{XW}{WX}\|_F^2 = (Y - XW)^\top (Y - XW)$$

$$= \text{tr}((Y - XW)^\top (Y - XW))$$

$$= \text{tr}(Y^\top Y - Y^\top XW - W^\top X^\top Y + W^\top X^\top XW)$$

$$\bullet dL = -\text{tr}(Y^\top X dW) + \underbrace{dW^\top X^\top Y}_{\cancel{Y^\top X dW}} + d\text{tr}(W^\top X^\top XW)$$

$$= -2\text{tr}(Y^\top X dW) + \frac{\text{tr}(dW^\top X^\top XW + W^\top X^\top X dW)}{2\text{tr}(W^\top X^\top X dW)}$$

$$Y^\top X = W^\top X^\top X$$

$$Y^\top X (X^\top X)^{-1} = W^\top$$

# Identification theorems

- Sometimes useful to go back and forth between differential & ordinary notations

► not always possible: e.g.,  $d(X^T X) = dX^T X + X^T dX$

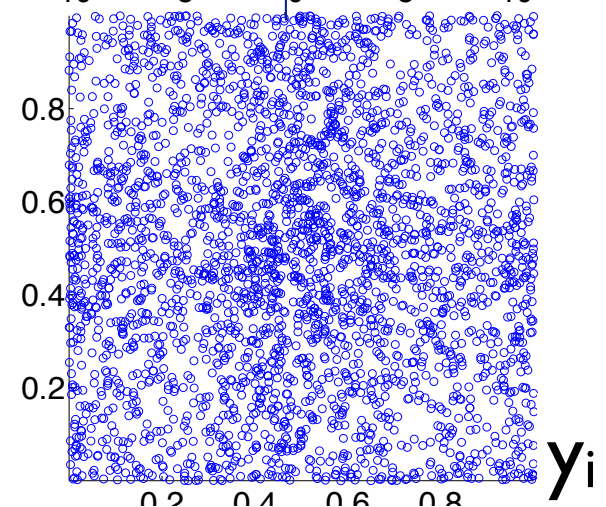
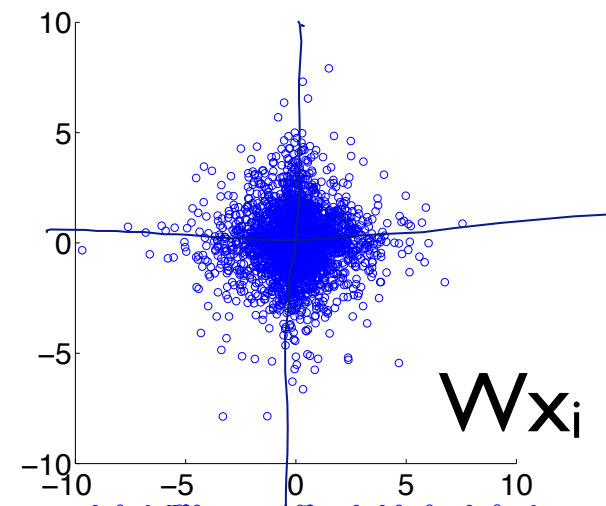
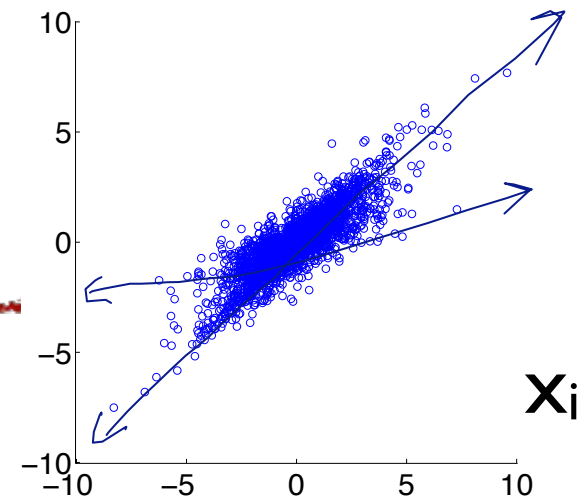
- Six common cases (ID thms):

ID for $df(x)$	scalar $x$	vector $\mathbf{x}$	matrix $X$
scalar $f$	$df = \textcircled{a} dx$	$df = \mathbf{a}^T d\mathbf{x}$	$df = \text{tr}(A^T dX)$
vector $\mathbf{f}$	$d\mathbf{f} = \mathbf{a} dx$	$d\mathbf{f} = A d\mathbf{x}$	$\textcircled{\phantom{0}}$
matrix $F$	$dF = A dx$	$\textcircled{\phantom{0}}$	$\textcircled{\phantom{0}}$

$df/dx$  → velocity  
 $\nabla$  → gradient  
 $\textcircled{\phantom{0}}$  → ~~Jacobian~~ Jacobian

# Ex: Infomax ICA

- Training examples  $\mathbf{x}_i \in \mathbb{R}^d, i = 1:n$
- Transformation  $\mathbf{y}_i = g(\mathbf{W}\mathbf{x}_i)$ 
  - ▶  $\mathbf{W} \in \mathbb{R}^{d \times d}$  *parameter*
  - ▶  $g(\mathbf{z}) =$  *scalar fn, componentwise*
- Want: *independent components*



# Ex: Infomax ICA

- $y_i = g(Wx_i)$

- ▶  $dy_i = J(x_i, W) dx_i = J_i dx_i$

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

- ▶ where  $P(y_i) = P(x_i) / |\det J(x_i, W)|$

$$\max_W \sum_i (\ln |\det J(x_i, W)| - \ln P(x_i))$$

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

