

Exploiting Feature Covariance in High-Dimensional Online Learning

Justin Ma (UCSD → Berkeley), Alex Kulesza (UPenn),

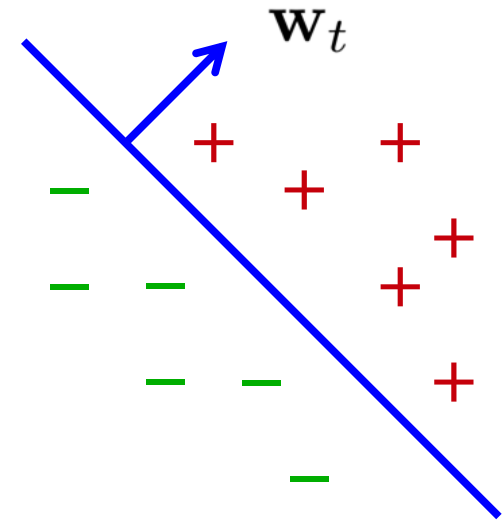
Mark Dredze (Johns Hopkins), Koby Crammer (The Technion),
Lawrence K. Saul (UCSD), Fernando Pereira (Google)

Presentation for AISTATS

May 14, 2010

Online Learning of Linear Classifiers

- Input \mathbf{x}_t
- Predict $\hat{y}_t = \text{sign}(\mathbf{w}_t \cdot \mathbf{x}_t)$
- Receive label $y_t \in \{-1, +1\}$
- Record error if $y_t \neq \hat{y}_t$
- Modify \mathbf{w}_t



Real-World Motivations

- Industrial-scale applications
 - **Large** data sets
($\sim 10^6$ — 10^9 examples)
 - **Nonstationary** data,
drifting concepts
- Attractions of online learning
 - **Single** pass over data
 - **Incremental** update
 - **Low** overhead in storage & compute



High-Dimensional Applications

- Ex: sentiment classification, malicious URL detection, Web spam
- Bag-of-words $\sim 10^6$ features
- **New features** introduced over time



Sentiment classification



Malicious URLs

Which online algorithm?

Perceptron [Rosenblatt, 1958]

Stochastic gradient [generalization in Bottou, 1998]

Bayesian logistic regression [MacKay, 1992] [Jaakkola & Jordan, 2000]

Online convex programming [Zinkevich, 2003]

Second-order perceptron [Cesa-Bianchi et al, 2005]

Passive-aggressive [Crammer et al, 2006]

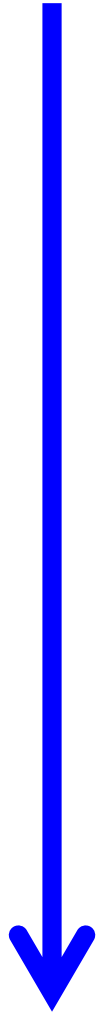
Confidence-weighted [Dredze et al, 2008]

Online ellipsoid method [Yang et al, 2009]

AROW [Crammer et al, 2009]

AdaGrad [Duchi et al, 2010]

Trend toward more complex updates
...(e.g., using **2nd-order** information)



Which online algorithm?

Perceptron [Rosenblatt, 1958]

Stochastic gradient [generalization in Bottou, 1998]

Bayesian logistic regression [MacKay, 1992] [Jaakkola & Jordan, 2000]

Online convex programming [Zinkevich, 2003]

Second-order perceptron [Cesa-Bianchi et al, 2005]

Passive-aggressive [Crammer et al, 2006]

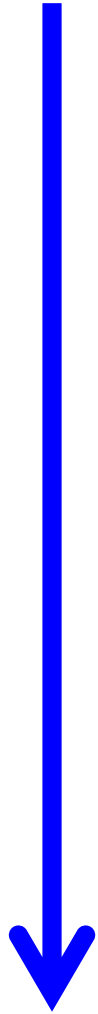
Confidence-weighted [Dredze et al, 2008]

Online ellipsoid method [Yang et al, 2009]

AROW [Crammer et al, 2009]

AdaGrad [Duchi et al, 2010]

Trend toward more complex updates
...(e.g., using **2nd-order** information)

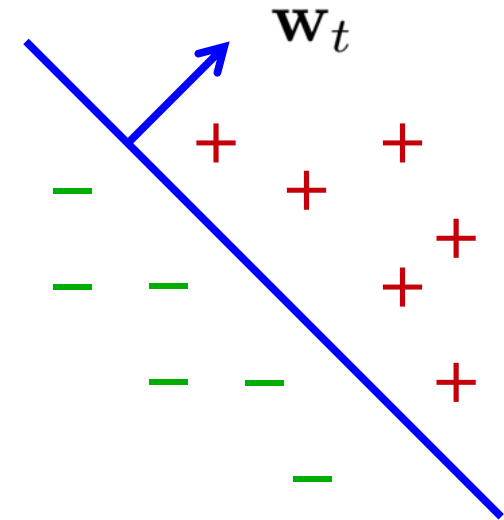


Perceptron

[Rosenblatt, 1958]

.Per-mistake update:

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_t \mathbf{X}_t$$



.Convergence in finite rounds for separable data

Passive-Aggressive (PA) Algorithm

[Crammer et al., 2006]

.Constrained optimization

$$\mathbf{w}_{t+1} \leftarrow \underset{\mathbf{w}}{\operatorname{argmin}} \quad \frac{1}{2} \|\mathbf{w}_t - \mathbf{w}\|^2$$
$$\text{s.t.} \quad y_t(\mathbf{w} \cdot \mathbf{x}_t) \geq 1$$

.Closed-form update

Amount of error

Proportional

$$\alpha_t = \max \left\{ \frac{1 - y_t(\mathbf{w}_t \cdot \mathbf{x}_t)}{\|\mathbf{x}_t\|^2}, 0 \right\}$$

$$\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + \alpha_t y_t \mathbf{x}_t$$

Confidence-Weighted (CW) Learning

[Dredze et al., 2008] [Crammer et al., 2009]

.Gaussian distribution over weight vector:

$$\mathbf{w}_t \sim \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t^{-1})$$

How to represent this matrix?

.Constrained problem:

$$\begin{aligned} (\boldsymbol{\mu}_{t+1}, \boldsymbol{\Sigma}_{t+1}^{-1}) &\leftarrow \underset{\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}}{\operatorname{argmin}} \operatorname{KL}(\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}^{-1}) \parallel \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t^{-1})) \\ \text{s.t. } &\Pr[y_t(\mathbf{w} \cdot \mathbf{x}_t) \geq 0] \geq \eta \end{aligned}$$

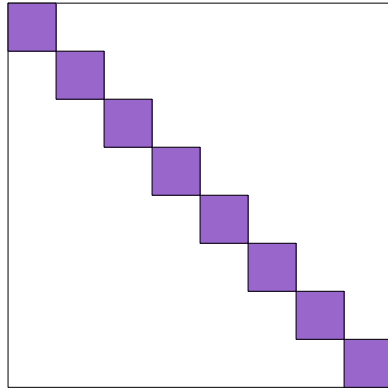
.Closed-form update:

$$\begin{aligned} \boldsymbol{\mu}_{t+1} &\leftarrow \boldsymbol{\mu}_t + \alpha_t y_t \boldsymbol{\Sigma}_t \mathbf{x}_t \\ \boldsymbol{\Sigma}_{t+1}^{-1} &\leftarrow \boldsymbol{\Sigma}_t^{-1} + \frac{\alpha_t \phi}{\sqrt{u_t}} \mathbf{x}_t \mathbf{x}_t^\top \end{aligned}$$

Update features at different rates

Representing Correlations

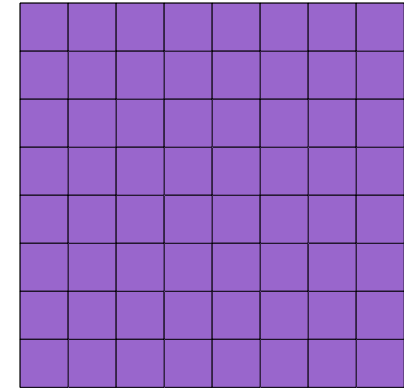
$$\Sigma^{-1} =$$



Diagonal

- ✓ $O(n)$ storage
- ✓ Low compute time

$$\Sigma^{-1} =$$



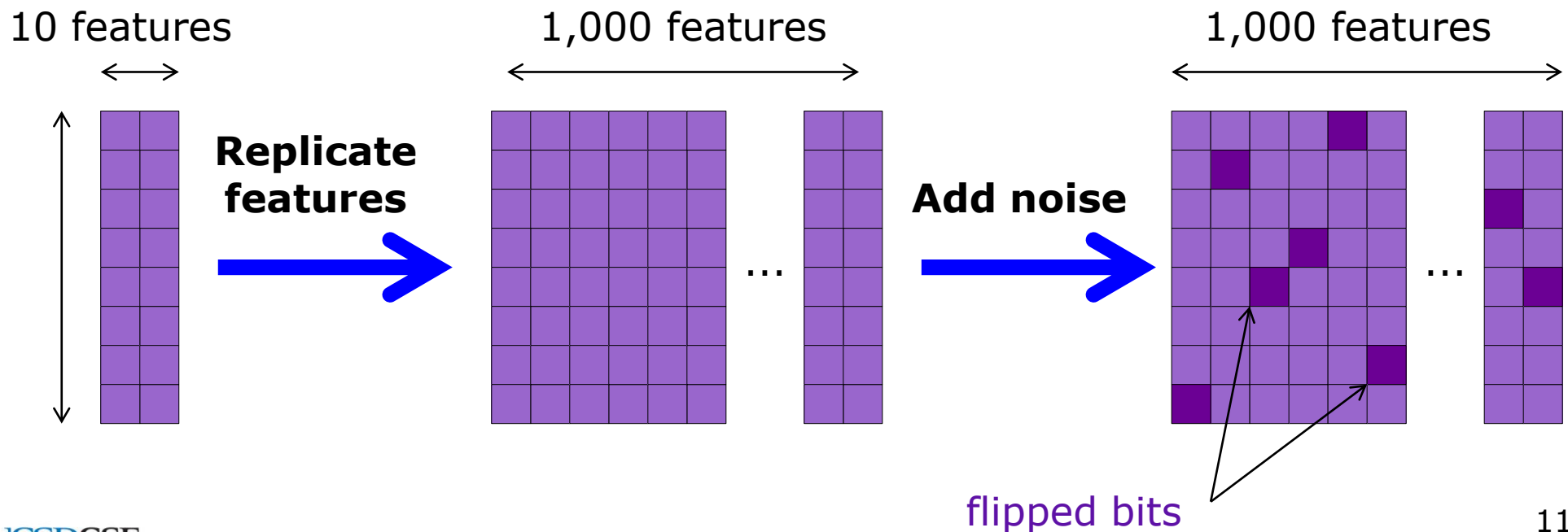
Full

- ✗ $O(n^2)$ storage
- ✗ High compute time

Why bother with full?

Benefits of Full

- When do we benefit from **full** covariance?
- Synthetic experiment: **noisy correlated** features
 - 100 runs, 1,000 examples, 1,000 binary features
 - 5% of features flipped



Noisy Correlated Features



Perceptron

PA

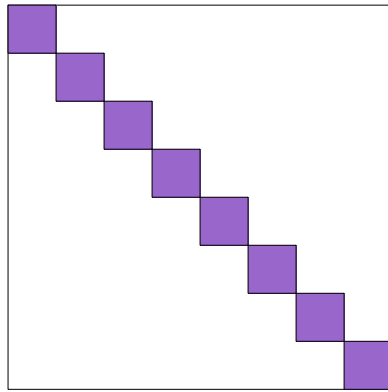
CW-diag

CW-full

- **Full** exploits correlations

Representing Correlations

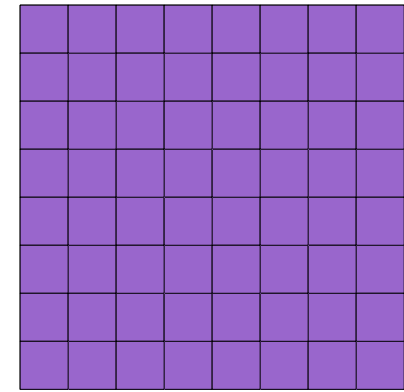
$$\Sigma^{-1} =$$



Diagonal

- ✓ $O(n)$ storage
- ✓ Low compute time
- ✓ Ignores correlations

$$\Sigma^{-1} =$$



Full

- ✗ $O(n^2)$ storage
- ✗ High compute time
- ✗ Exploits correlations

Can obtain benefits of both?

Factored Approximation

$$\Sigma^{-1} = \mathbf{D} + \mathbf{R}\mathbf{R}^T$$

The diagram shows the equation $\Sigma^{-1} = \mathbf{D} + \mathbf{R}\mathbf{R}^T$. On the left, a blue diagonal matrix \mathbf{D} is shown with a width of n and a height of n . To its right is a plus sign. Further right is a red matrix \mathbf{R} with a width of n and a height of k . To its right is another red matrix \mathbf{R}^T with a width of n and a height of k . The overall structure represents the decomposition of the inverse covariance matrix into a diagonal matrix and a low-rank matrix product.

- Approx. full inv. covariance (called precision)
- $O(kn)$ storage (k = number of factors)
- Compress matrix updates – **factor analysis**

Factor Analysis

- Exact and approximate distributions

$$\rightarrow P(\mathbf{w}_t) = \mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\Sigma}_t^{-1})$$

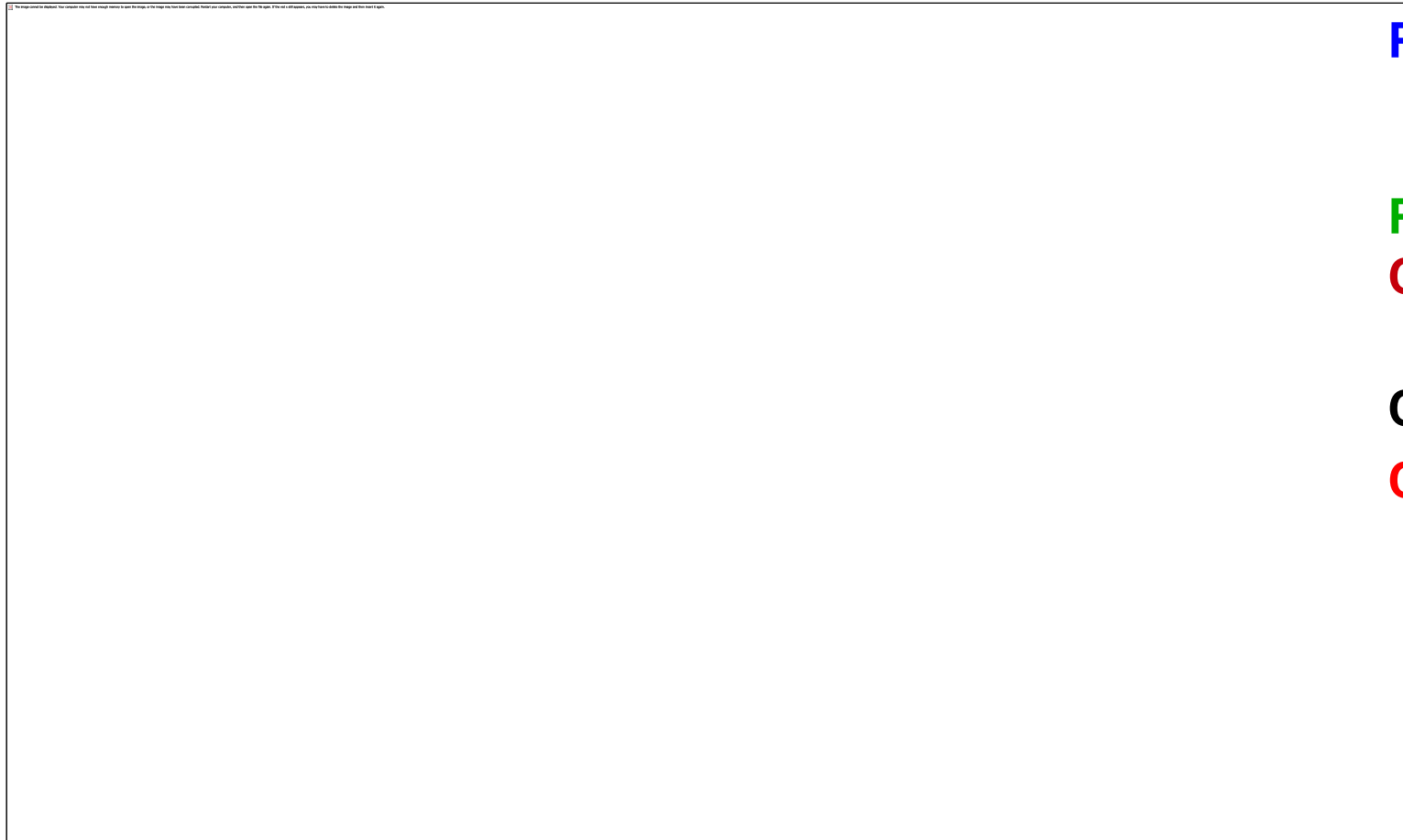
$$\rightarrow \hat{P}(\mathbf{w}_t) = \mathcal{N}(\boldsymbol{\mu}_t, \mathbf{D}_t + \mathbf{R}_t \mathbf{R}_t^\top)$$

- Minimize KL divergence using EM procedure

$$\min \text{KL}(P_t, \hat{P}_t)$$

- Computation cost: $O(nk^2)$

Synthetic with CW-fact



Perceptron

PA

CW-diag

CW-fact4

CW-full

- **Full** exploits correlations
- **Factored** approximates full
- **40x** less memory

Benefits of **approximating** full in real-world applications?

Detecting Malicious URLs

- Live feature collection of URLs
- Per trial: 200,000 examples, 10^6 features (mostly binary)

**Additional mistakes
by CW-diag**

Perceptron

PA

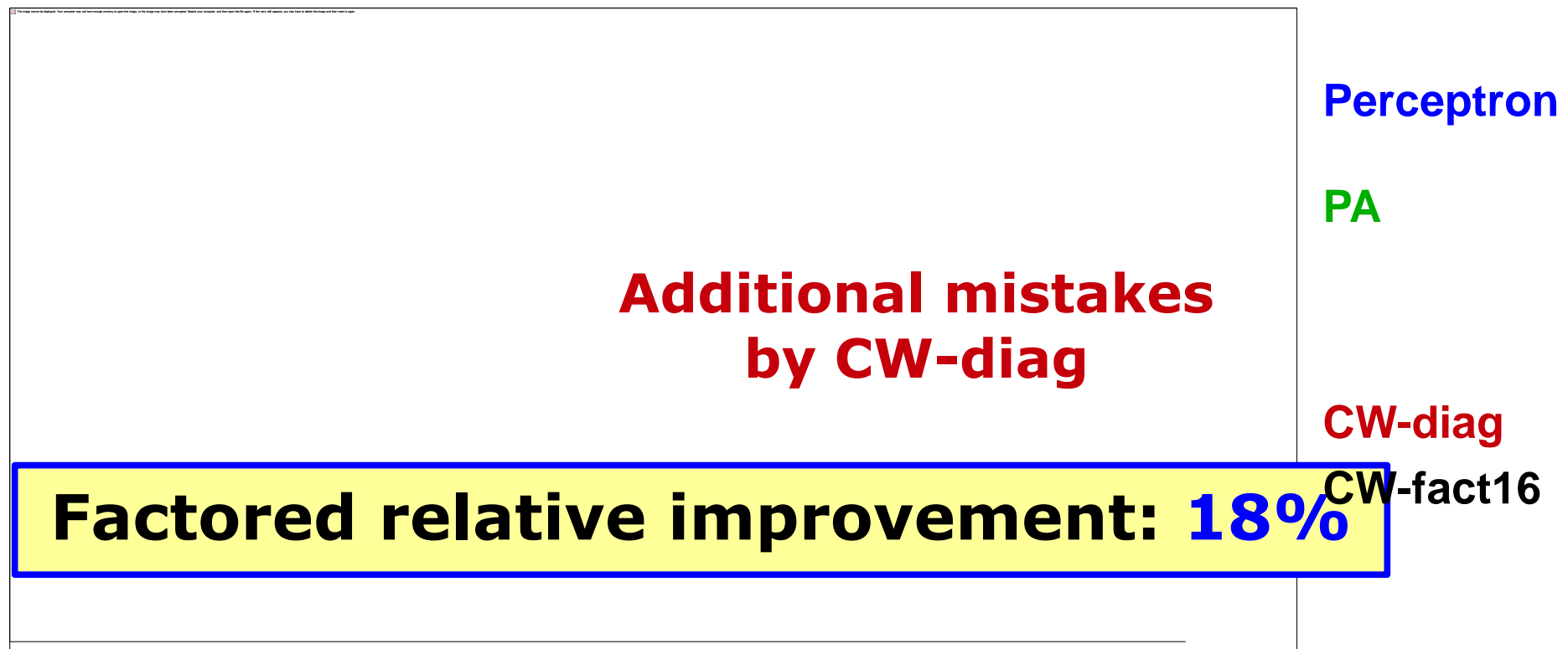
CW-diag

CW-fact16

Factored relative improvement: 5%

Web Spam Classification

- Web pages [PASCAL Large-Scale Learning competition]
- Per-run: 175,000 examples, 680,000 features (text 3-grams)



Document Classification

[Dredze et al., 2008]

- 2,000–18,000 examples, 10^4 – 10^6 features

**Improvement by
CW-fact8**

Reuters

20 Newsgroups

Sentiment

Conclusion

- Full and factored covariance help when...
 - features are correlated
 - $\# \text{features} > \# \text{examples}$
(not needed when $\# \text{examples} > \# \text{features}$ [Sec. 2])
- Factored improves high-dimensional apps
 - NLP, URL classification, Web spam, others
- Future work
 - Correlation modeling in other online algorithms?
 - Other ways to model correlation structure?

Code is Available

<http://sysnet.ucsd.edu/projects/url/>

- Coauthors in Sardinia...



Alex

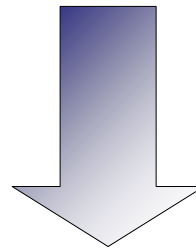


Koby

Update Compression

$$\Sigma^{-1} = \begin{matrix} \text{[Blue diagonal matrix]} \\ \mathbf{D} \end{matrix} + \begin{matrix} \text{[Red matrix]} \\ \mathbf{RR}^T \end{matrix} + \begin{matrix} \text{[Purple matrix]} \\ \delta \mathbf{xx}^T \end{matrix}$$

Factor Analysis



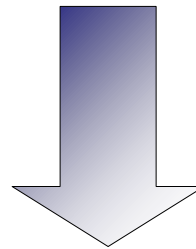
$$\Sigma^{-1} = \begin{matrix} \text{[Magenta diagonal matrix]} \\ \mathbf{D} \end{matrix} + \begin{matrix} \text{[Magenta matrix]} \\ \mathbf{RR}^T \end{matrix}$$

Buffering Updates

$$\Sigma^{-1} = \begin{matrix} \text{[Blue Diagonal Matrix]} & + & \begin{matrix} \text{[Red Column]} & \text{[Red Row]} \end{matrix} & + & \begin{matrix} \text{[Purple Column]} & \text{[Purple Row]} \end{matrix} \end{matrix}$$

D **RR^T** **BB^T**

Factor Analysis



$$\Sigma^{-1} = \begin{matrix} \text{[Magenta Diagonal Matrix]} & + & \begin{matrix} \text{[Magenta Column]} & \text{[Magenta Row]} \end{matrix} \end{matrix}$$

D **RR^T**