

# Learning with Marginalized Corrupted Features

---

Laurens van der Maaten

Minmin Chen

Stephen Tyree

Kilian Weinberger

# Classification of text and image data

---

**Classify....**

# Classification of text and image data

## Classify....

### Man Googles Matt Damon's Address Because, Well, He's Crazy And Wants To Murder Him

NOVEMBER 8, 2012 | ISSUE 48-45 | [MORE NEWS](#)



The 29-year-old Easter, who stared at this picture of Matt Damon for two hours because, well, he's mentally ill.

SALISBURY, MARYLAND—After rereading actor Matt Damon's Wikipedia page for the 13th time since 9 a.m. today, local man Dan Easter decided to look up the celebrity's home address on Google because, well, he's admittedly crazy and wants to murder him.

Saying he planned to "just click around" a couple websites to see if the *Bourne Identity* star's address was listed anywhere on the Internet, Easter told reporters that, you know, he's ultimately a mentally ill madman who wants to break into Matt Damon's house in the dead of night and, you guessed it, kill him in front of his wife and children.

"I figured I would just type Matt Damon's name into Google because, to make a long story short, I'm psychologically disturbed and I want to assassinate him," said the 29-year-old man, who, by his own admission, is extremely unstable and has absolutely no business being anywhere other than a mental institution. "I see him in movies and magazines all the time, and it made me wonder where he lives. I'm also clinically insane. That's why I want to see these and threaten him to death."

ARTICLE TOOLS

Tweet 120

Like 954

+1 16

tumblr +

RELATED ARTICLES

[Fox Voluntarily Removes Reality From Programming](#)

[TV Viewers Outraged At Timing Of Commercial Break](#)

... documents  
by topic

# Classification of text and image data

## Classify....

### Man Googles Matt Damon's Address Because, Well, He's Crazy And Wants To Murder Him

NOVEMBER 8, 2012 | ISSUE 48-45 | [MORE NEWS](#)



... documents  
by topic

### Kindle vs. Nook vs. iPad: Which e-book reader should you buy?

With ultraaffordable e-ink readers, midprice color tablets like the Nexus 7 and Kindle Fire, and even the more expensive iPads all vying for your e-book dollar, what's the best choice for you? It depends.



by John P. Falcone | August 7, 2012 1:08 PM PDT

[Follow](#)



6.3K



1.5K



200



+1



282



More +

Comments

573

The 29-year-old Easter, who stared at this picture of Matt Damon for two hours before SALISBURY, MARYLAND—After rereading actor Matt Damon's page for the 13th time since 9 a.m. today, local man Dan Easter look up the celebrity's home address on Google because, well admittedly crazy and wants to murder him.

Saying he planned to "just click around" a couple websites to see if Matt Damon's address was listed anywhere on the Internet, Easter told reporters that, you know, he's ultimately a mentally ill mad man who wants to break into Matt Damon's house in the dead of night and kill him in front of his wife and children.

"I figured I would just type Matt Damon's name into Google and see what comes up. Well, I'm psychologically disturbed and I want to assassinate him," said the 29-year-old man, who, by his own admission, is extremely unstable and has absolutely no business being anywhere other than a mental institution. "I see him in movies and magazines and it made me wonder where he lives. I'm also clinically insane and I want to see those and strangle him to death."



(Credit: Sarah Tew/CNET)

... documents  
by sentiment

**Editors' note, September 7, 2012:** As of September 6, Amazon has announced [all-new Kindle e-readers and tablets](#) for 2012 that dramatically offer the buying decisions listed below. The first wave of the new Amazon products are due to ship by September 14. We'll update this story in detail after we review those models. By that time, we'll also find out what Apple is announcing at its [September 12 event](#). In the meantime, the competition includes the [Kindle Fire](#) and the [Nook HD](#).

# Classification of text and image data

## Classify....

### Man Googles Matt Damon's Address Because, Well, He's Crazy And Wants To Murder Him

NOVEMBER 8, 2012 | ISSUE 48-45 | MORE NEWS



... documents  
by topic

### Kindle vs. Nook vs. iPad: Which e-book reader should you buy?

With ultraaffordable e-ink readers, midprice color tablets like the Nexus 7 and Kindle Fire, and even the more expensive iPads all vying for your e-book dollar, what's the best choice for you? It depends.



by John P. Falcone | August 7, 2012 1:08 PM PDT

Follow

6.3K 1.5K 200 +1 282 More +

Comments 573



(Credit: Sarah Tew/CNET)

**Editors' note, September 7, 2012:** As of September 6, Amazon has announced [all-new Kindle e-readers and tablets](#) for 2012 that dramatically offer the buying decisions listed below. The first wave of the new Amazon products are due to ship by September 14. We'll update this story in detail after we review those models. By that time, we'll also find out what Apple is announcing at its [September 12 event](#). In the meantime, the competition includes [Kindle](#), [Nook](#), [iPads](#), [Kindle Fire](#), [Nexus 7](#), and [Kindle](#).



... images  
by object

... documents  
by sentiment

# Empirical risk minimization

---

- Learn model based on annotated data  $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$  by minimizing:

$$\min_{\Theta} \mathcal{L}(\mathcal{D}; \Theta) = \sum_{n=1}^N L(\mathbf{x}_n, y_n; \Theta)$$

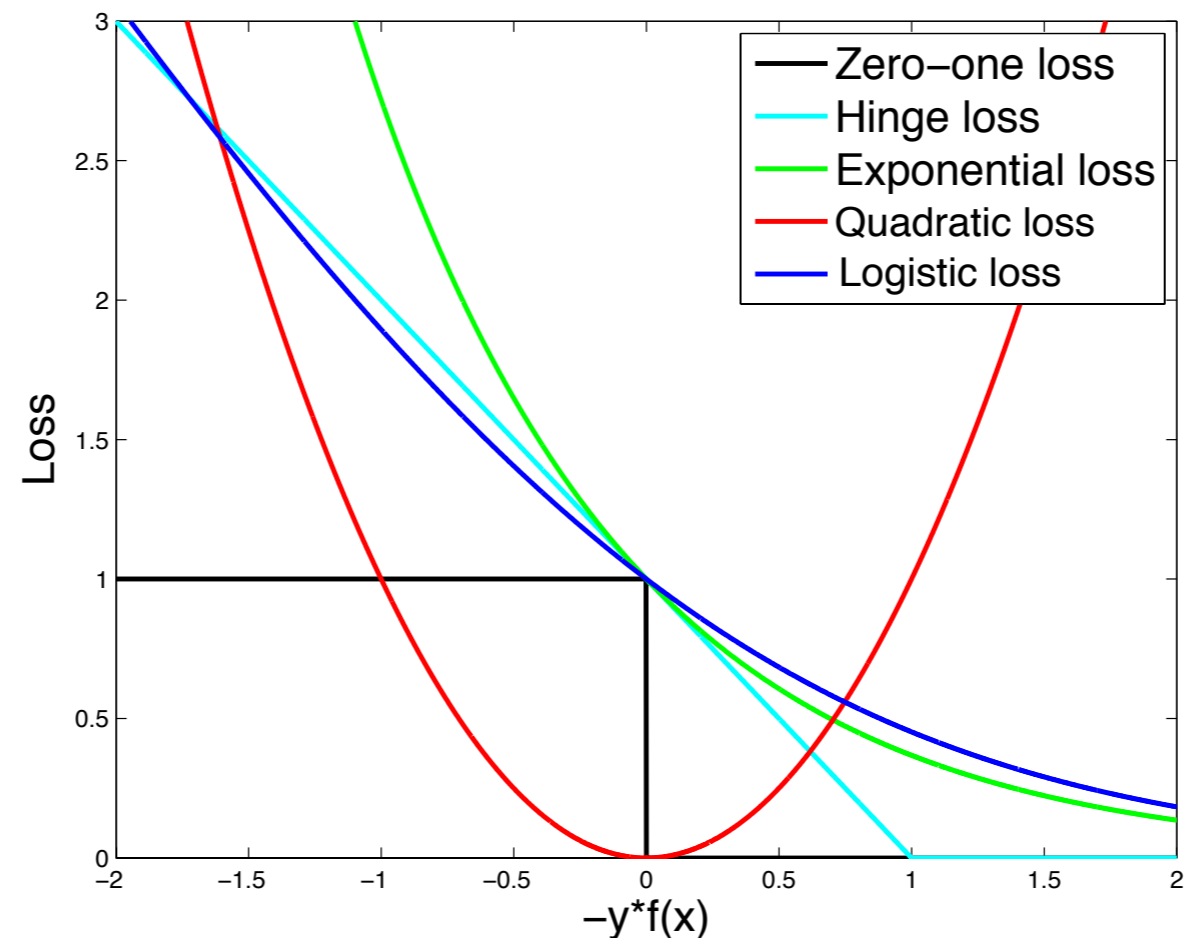
# Empirical risk minimization

---

- Learn model based on annotated data  $\mathcal{D} = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$  by minimizing:

$$\min_{\Theta} \mathcal{L}(\mathcal{D}; \Theta) = \sum_{n=1}^N L(\mathbf{x}_n, y_n; \Theta)$$

- Herein, typical examples of the loss function include:



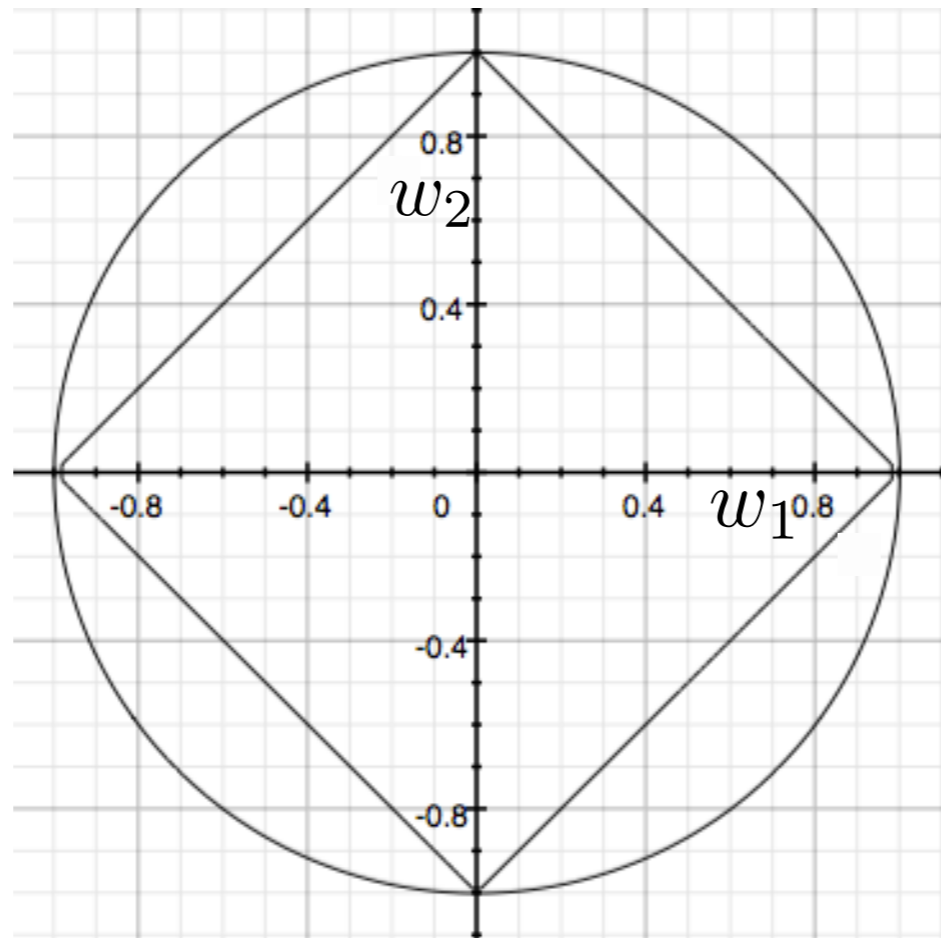
# Regularization by priors

---

- *Regularizers* incorporate a term in the loss that penalizes complex models:

$$\tilde{\mathcal{L}}(\mathbf{x}, y; \mathbf{w}) = \mathcal{L}(\mathbf{x}, y; \mathbf{w}) + \lambda \mathcal{R}(\mathbf{w})$$

e.g.,  $\mathcal{R}(\mathbf{w}) = \|\mathbf{w}\|^2$  or  $\mathcal{R}(\mathbf{w}) = |\mathbf{w}|$





# Regularization by priors

---

- Getting the right regularizer is *tricky!*

# Regularization by priors

---

- Getting the right regularizer is *tricky!*
  - Most norm-based regularizers are rather *arbitrary*:
    - L1 and L2-regularization are popular mainly for computational reasons

# Regularization by priors

---

- Getting the right regularizer is *tricky!*
  - Most norm-based regularizers are rather *arbitrary*:
    - L1 and L2-regularization are popular mainly for computational reasons
- Most practitioners have *bad intuitions* about model parameters...
  - ... but they do understand their data!

# Regularization by priors

---

- Getting the right regularizer is *tricky!*

Instead of restricting the  
**model parameters**,  
can't we incorporate  
**knowledge about the data**  
instead?

- Most normal

- L1 and L2

- Most practical

- ... but they do understand their data!

# Movie reviews

- Are these reviews positive or negative?



This is a boring movie with a lot of decadence and bad influence on people. I can't believe this movie won awards! I would not recommend this though it's so famous.



The movie is great, and in perfect condition. Came in time. I'd recommend the movie itself, and I would purchase movies from here again.



This movie is awesome, if you have not seen Tarrantino movies on Blu Ray you are missing out. Blu Ray brings these movies to life, especially if you have a good surround sound system.



I tried to watch. I bought it because of the Micah quote. If you like to watch people get high and talk filthy this is for you.

The screenshot shows the Amazon.com product page for the movie "Pulp Fiction". The browser address bar shows the URL: [www.amazon.com/Pulp-Fiction/dp/B005T3AYAE/ref=sr\\_1\\_1?ie=UTF8&qid=1364849860](http://www.amazon.com/Pulp-Fiction/dp/B005T3AYAE/ref=sr_1_1?ie=UTF8&qid=1364849860). The page features the Amazon logo, navigation links like "Join Prime", "Today's Deals", "Gift Cards", "Sell", and "Help". The search bar contains "pulp fiction". The product title is "Pulp Fiction" with a rating of 4.5 stars from 1,089 reviews. The description states: "Writer/director Quentin Tarantino delivers an unforgettable cast of characters -- including a pair of low-rent hit men, their boss's sexy wife, and a desperate prizefighter -- in a wildly entertaining and exhilarating motion picture adventure that both thrills and amuses!". The starring cast includes John Travolta, Samuel L. Jackson, Uma Thurman, and Harvey Keitel. The director is Quentin Tarantino, the runtime is 2 hours 35 minutes, and the release year is 1994. The studio is Lionsgate. On the right, there are options for "24 hour rental" for \$1.99 and "Buy movie" for \$9.99. There is also an "Add to Watchlist" button and a "Play trailer" button at the bottom.

# Regularization by corruption

- Remove each word with probability  $q$ :



This is a boring movie with a lot of decadence and bad influence on people. I can't believe this movie won awards! I would not recommend this though it's so famous.



This is a boring [REDACTED] with a lot of decadence and bad influence on people. I [REDACTED] believe this [REDACTED] won awards! I would not recommend this though it's so [REDACTED].



This is a [REDACTED] movie with a lot of decadence and bad [REDACTED] on [REDACTED]. I can't believe this movie won awards! I would not [REDACTED] this though it's so famous.



The movie is great, and in perfect condition. Came in time. I'd recommend the movie itself, and I would purchase movies from here again.



The [REDACTED] is great, and in [REDACTED] condition. Came in time. I'd recommend the [REDACTED] itself, and I would [REDACTED] movies from here again.



The movie is [REDACTED], and in perfect condition. Came in time. I'd [REDACTED] the movie itself, and I would [REDACTED] movies from here again.



This movie is awesome, if you have not seen Tarrantino movies on Blu Ray you are missing out. Blu Ray brings these movies to life, especially if you have a good surround sound system.



This movie is [REDACTED] if you [REDACTED] not seen Tarrantino movies on Blu Ray you are missing out. Blu Ray brings these [REDACTED] to life, especially if [REDACTED] have a good [REDACTED] sound system.



This movie is [REDACTED], if you have [REDACTED] seen Tarrantino movies on Blu Ray you are missing out. Blu Ray brings these movies to life, [REDACTED] if you have a good surround sound system.



I tried to watch. I bought it because of the Micah quote. If you like to watch people get high and talk filthy this is for you.



I tried to [REDACTED]. I bought it [REDACTED] of the [REDACTED] quote. If you like to watch people get high and talk [REDACTED] this is for you.



I tried to watch. I [REDACTED] it because of the Micah quote. If you like to watch people get high an [REDACTED] filthy this is for [REDACTED]

# Regularization by corruption

---

- Define *label-invariant corruptions* that can be applied to the data
- Training on such corrupted data leads to *robustness* to the corruption
- *Robustness* is intimately related to *regularization* of the model
- We show that this can be done *efficiently* by *marginalizing* over corruptions

# Regularization by corruption

---

- Instead of regularizer, define a *label-invariant corrupting distribution*:

This is a boring [redacted] with a lot of decadence and bad influence on people. I [redacted] believe this [redacted] won awards! I would not recommend this though it's so [redacted].

This is a boring movie with a lot of decadence and bad influence on people. I can't believe this movie won awards! I would not recommend this though it's so famous.

$$p(\tilde{\mathbf{x}}|\mathbf{x}) = \prod_{d=1}^D p(\tilde{x}_d|x_d; \eta_d), \text{ with } \mathbb{E}[\tilde{\mathbf{x}}]_{p(\tilde{\mathbf{x}}|\mathbf{x})} = \mathbf{x}$$

- We will assume the corruption are independent across features (this assumption may be relaxed for Gaussian corruptions)



# Regularization by corruption

- Instead of regularizer, define a *label-invariant corrupting distribution*:

This is a boring [redacted] with a lot of decadence and bad influence on people. I [redacted] believe this [redacted] won awards! I would not recommend this though it's so [redacted].

This is a boring movie with a lot of decadence and bad influence on people. I can't believe this movie won awards! I would not recommend this though it's so famous.

$$p(\tilde{\mathbf{x}}|\mathbf{x}) = \prod_{d=1}^D p(\tilde{x}_d|x_d; \eta_d), \text{ with } \mathbb{E}[\tilde{\mathbf{x}}]_{p(\tilde{\mathbf{x}}|\mathbf{x})} = \mathbf{x}$$

- We will assume the corruption are independent across features (this assumption may be relaxed for Gaussian corruptions)

Distribution	PDF
Blankout noise	$p(\tilde{x}_{nd} = 0) = q_d$ $p(\tilde{x}_{nd} = \frac{1}{1-q_d} x_{nd}) = 1 - q_d$
Gaussian noise	$p(\tilde{x}_{nd} x_{nd}) = \mathcal{N}(\tilde{x}_{nd} x_{nd}, \sigma^2)$
Laplace noise	$p(\tilde{x}_{nd} x_{nd}) = Lap(\tilde{x}_{nd} x_{nd}, \lambda)$
Poisson noise	$p(\tilde{x}_{nd} x_{nd}) = Poisson(\tilde{x}_{nd} x_{nd})$

This is a boring [redacted] with a lot of decadence and bad influence on people. I [redacted] believe this [redacted] won awards! I would not recommend this though it's so [redacted].



keep	0
amazing	<del>7</del>
ideas	2
value	0
poor	0
...	...
average	1

5

# Simple approach

---

- For each example, generate  $M$  corrupted examples and use these as data
- This amounts to minimizing the loss on an *augmented, corrupted* training set:

$$\mathcal{L}(\tilde{\mathcal{D}}; \Theta) = \sum_{n=1}^N \frac{1}{M} \sum_{m=1}^M L(\tilde{\mathbf{x}}_{nm}, y_n; \Theta) \text{ with } \tilde{\mathbf{x}}_{nm} \sim p(\tilde{\mathbf{x}}_n | \mathbf{x}_n)$$

This is a [REDACTED] movie with a lot of  
decadence and bad [REDACTED] on

[REDACTED]. I can  
awards! I would not  
though it's so [REDACTED].

This is a boring [REDACTED] with a lot of  
decadence and bad influence on  
people. I [REDACTED] believe this [REDACTED] won  
awards! I would not recommend this  
though it's so [REDACTED].

...

# Simple approach

---

- For each example, generate  $M$  corrupted examples and use these as data\*
- This amounts to minimizing the loss on an *augmented, corrupted* training set:

$$\mathcal{L}(\tilde{\mathcal{D}}; \Theta) = \sum_{n=1}^N \frac{1}{M} \sum_{m=1}^M L(\tilde{\mathbf{x}}_{nm}, y_n; \Theta) \text{ with } \tilde{\mathbf{x}}_{nm} \sim p(\tilde{\mathbf{x}}_n | \mathbf{x}_n)$$

- This quickly gets computationally prohibitive, unless...

# Marginalized Corrupted Features

---

- For each example, generate  $M$  corrupted examples and use these as data
- This amounts to minimizing the loss on an *augmented, corrupted* training set:

$$\mathcal{L}(\tilde{\mathcal{D}}; \Theta) = \sum_{n=1}^N \frac{1}{M} \sum_{m=1}^M L(\tilde{\mathbf{x}}_{nm}, y_n; \Theta) \text{ with } \tilde{\mathbf{x}}_{nm} \sim p(\tilde{\mathbf{x}}_n | \mathbf{x}_n)$$

- This quickly gets computationally prohibitive, unless  $M \rightarrow \infty$
- Law of large numbers leads to the *expected loss under the corruption model*:

$$\mathcal{L}(\mathcal{D}; \Theta) = \sum_{n=1}^N \mathbb{E}[L(\tilde{\mathbf{x}}_n, y_n; \Theta)]_{p(\tilde{\mathbf{x}}_n | \mathbf{x}_n)}$$

# Quadratic loss

---

- Working out the MCF expectation (for independent corruption) gives:

$$\begin{aligned}\mathcal{L}(\mathcal{D}; \mathbf{w}) &= \sum_{n=1}^N \mathbb{E} \left[ (\mathbf{w}^T \tilde{\mathbf{x}}_n - y_n)^2 \right]_{p(\tilde{\mathbf{x}}_n | \mathbf{x}_n)} \\ &= \mathbf{w}^T \left( \sum_{n=1}^N \mathbb{E}[\tilde{\mathbf{x}}_n] \mathbb{E}[\tilde{\mathbf{x}}_n]^T + V[\tilde{\mathbf{x}}_n] \right) \mathbf{w} - 2 \left( \sum_{n=1}^N y_n \mathbb{E}[\tilde{\mathbf{x}}_n] \right)^T \mathbf{w} + N\end{aligned}$$

- Practical if we can compute the *mean* and *variance* of corrupting distribution

# Quadratic loss

---

- Working out the MCF expectation (for independent corruption) gives:

$$\begin{aligned}\mathcal{L}(\mathcal{D}; \mathbf{w}) &= \sum_{n=1}^N \mathbb{E} \left[ (\mathbf{w}^T \tilde{\mathbf{x}}_n - y_n)^2 \right]_{p(\tilde{\mathbf{x}}_n | \mathbf{x}_n)} \\ &= \mathbf{w}^T \left( \sum_{n=1}^N \mathbb{E}[\tilde{\mathbf{x}}_n] \mathbb{E}[\tilde{\mathbf{x}}_n]^T + V[\tilde{\mathbf{x}}_n] \right) \mathbf{w} - 2 \left( \sum_{n=1}^N y_n \mathbb{E}[\tilde{\mathbf{x}}_n] \right)^T \mathbf{w} + N\end{aligned}$$

- Practical if we can compute the *mean* and *variance* of corrupting distribution
- The objective function remains *convex*; optimal solution given by:

$$\mathbf{w}^* = \left( \sum_{n=1}^N \mathbb{E}[\tilde{\mathbf{x}}_n] \mathbb{E}[\tilde{\mathbf{x}}_n]^T + V[\tilde{\mathbf{x}}_n] \right)^{-1} \left( \sum_{n=1}^N y_n \mathbb{E}[\tilde{\mathbf{x}}_n] \right)$$

# Quadratic loss

---

- Examples of corrupting distributions of interest:

Distribution	PDF	$\mathbb{E}[\tilde{\mathbf{x}}_{nd}]_{p(\tilde{\mathbf{x}}_{nd} \mathbf{x}_{nd})}$	$\mathbf{V}[\tilde{\mathbf{x}}_{nd}]_{p(\tilde{\mathbf{x}}_{nd} \mathbf{x}_{nd})}$
Blankout noise	$p(\tilde{\mathbf{x}}_{nd} = 0) = q_d$ $p(\tilde{\mathbf{x}}_{nd} = \frac{1}{1-q_d}\mathbf{x}_{nd}) = 1 - q_d$	$\mathbf{x}_{nd}$	$\frac{1}{1-q_d}\mathbf{x}_{nd}^2$
Gaussian noise	$p(\tilde{\mathbf{x}}_{nd} \mathbf{x}_{nd}) = \mathcal{N}(\tilde{\mathbf{x}}_{nd} \mathbf{x}_{nd}, \sigma^2)$	$\mathbf{x}_{nd}$	$\sigma^2$
Laplace noise	$p(\tilde{\mathbf{x}}_{nd} \mathbf{x}_{nd}) = Lap(\tilde{\mathbf{x}}_{nd} \mathbf{x}_{nd}, \lambda)$	$\mathbf{x}_{nd}$	$2\lambda^2$
Poisson noise	$p(\tilde{\mathbf{x}}_{nd} \mathbf{x}_{nd}) = Poisson(\tilde{\mathbf{x}}_{nd} \mathbf{x}_{nd})$	$\mathbf{x}_{nd}$	$\mathbf{x}_{nd}$

- Using Gaussian corruptions leads to an interesting special case:

$$\mathcal{L}(\mathcal{D}; \mathbf{w}) = \mathbf{w}^T \left( \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{w} - 2 \left( \sum_{n=1}^N y_n \mathbf{x}_n \right)^T \mathbf{w} + \sigma^2 N \mathbf{w}^T \mathbf{w} + N$$

- Minimizing *MCF-Gaussian quadratic loss* leads to *ridge regression*!

# Exponential loss

---

- Working out the MCF expectation (for independent corruption) gives:

$$\begin{aligned}\mathcal{L}(\mathcal{D}; \mathbf{w}) &= \sum_{n=1}^N \mathbb{E} \left[ \exp \left( -y_n \mathbf{w}^T \tilde{\mathbf{x}}_n \right) \right]_{p(\tilde{\mathbf{x}}_n | \mathbf{x}_n)} \\ &= \sum_{n=1}^N \prod_{d=1}^D \mathbb{E} \left[ \exp \left( -y_n w_d \tilde{x}_{nd} \right) \right]_{p(\tilde{x}_{nd} | x_{nd})}\end{aligned}$$



# Exponential loss

---

- Working out the MCF expectation (for independent corruption) gives:

$$\begin{aligned}\mathcal{L}(\mathcal{D}; \mathbf{w}) &= \sum_{n=1}^N \mathbb{E} \left[ \exp \left( -y_n \mathbf{w}^T \tilde{\mathbf{x}}_n \right) \right]_{p(\tilde{\mathbf{x}}_n | \mathbf{x}_n)} \\ &= \sum_{n=1}^N \prod_{d=1}^D \mathbb{E} \left[ \exp \left( -y_n w_d \tilde{x}_{nd} \right) \right]_{p(\tilde{x}_{nd} | x_{nd})}\end{aligned}$$

- This can be recognized as a product of *moment-generating functions*:

$$M_x(t) = \mathbb{E}[\exp(tx)], t \in \mathbb{R}$$

# Moment-generating functions

Moment-generating function - Wikipedia, the free encyclopedia

en.wikipedia.org/wiki/Moment-generating\_function

Moment-generating function - Wikipedia, the free encyclopedia

Here are some examples of the moment generating function and the characteristic function for comparison. It can be seen that the characteristic function is a [Wick rotation](#) of the moment generating function  $M_X(t)$  when the latter exists.

Distribution	Moment-generating function $M_X(t)$	Characteristic function $\phi(t)$
<a href="#">Bernoulli</a> $P(X = 1) = p$	$1 - p + pe^t$	$1 - p + pe^{it}$
<a href="#">Geometric</a> $(1 - p)^{k-1} p$	$\frac{pe^t}{1 - (1 - p)e^t}$ for $t < -\ln(1 - p)$	$\frac{pe^{it}}{1 - (1 - p)e^{it}}$
<a href="#">Binomial</a> $B(n, p)$	$(1 - p + pe^t)^n$	$(1 - p + pe^{it})^n$
<a href="#">Poisson</a> $\text{Pois}(\lambda)$	$e^{\lambda(e^t - 1)}$	$e^{\lambda(e^{it} - 1)}$
<a href="#">Uniform (continuous)</a> $U(a, b)$	$\frac{e^{tb} - e^{ta}}{t(b - a)}$	$\frac{e^{itb} - e^{ita}}{it(b - a)}$
<a href="#">Uniform (discrete)</a> $U(a, b)$	$\frac{e^{at} - e^{(b+1)t}}{(b - a + 1)(1 - e^t)}$	$\frac{e^{ait} - e^{(b+1)it}}{(b - a + 1)(1 - e^{it})}$
<a href="#">Normal</a> $N(\mu, \sigma^2)$	$e^{t\mu + \frac{1}{2}\sigma^2 t^2}$	$e^{it\mu - \frac{1}{2}\sigma^2 t^2}$
<a href="#">Chi-squared</a> $\chi_k^2$	$(1 - 2t)^{-k/2}$	$(1 - 2it)^{-k/2}$
<a href="#">Gamma</a> $\Gamma(k, \theta)$	$(1 - t\theta)^{-k}$	$(1 - it\theta)^{-k}$
<a href="#">Exponential</a> $\text{Exp}(\lambda)$	$(1 - t\lambda^{-1})^{-1}$	$(1 - it\lambda^{-1})^{-1}$
<a href="#">Multivariate normal</a> $N(\mu, \Sigma)$	$e^{t^T \mu + \frac{1}{2} t^T \Sigma t}$	$e^{it^T \mu - \frac{1}{2} t^T \Sigma t}$
<a href="#">Degenerate</a> $\delta_a$	$e^{ta}$	$e^{ita}$
<a href="#">Laplace</a> $L(\mu, b)$	$\frac{e^{t\mu}}{1 - b^2 t^2}$	$\frac{e^{it\mu}}{1 + b^2 t^2}$
<a href="#">Negative Binomial</a> $\text{NB}(r, p)$	$\frac{((1 - p)e^t)^r}{(1 - pe^t)^r}$	$\frac{((1 - p)e^{it})^r}{(1 - pe^{it})^r}$

# Blankout: Ensemble interpretation

---

- MCF with blankout has an interesting interpretation as an *ensemble*

# Blankout: Ensemble interpretation

---

- MCF with blankout has an interesting interpretation as an *ensemble*
- Example for model with two input features:

$$\begin{aligned} \mathcal{L}(\mathcal{D}; \mathbf{w}) = \sum_{n=1}^N & \left[ q_1 q_2 + (1 - q_1) q_2 \exp(-y_n w_1 x_{n1}) \right. \\ & + (1 - q_2) q_1 \exp(-y_n w_2 x_{n2}) \\ & \left. + (1 - q_1)(1 - q_2) \exp(-y_n [w_1 x_{n1} + w_2 x_{n2}]) \right] \end{aligned}$$

# Blankout: Ensemble interpretation

---

- MCF with blankout has an interesting interpretation as an *ensemble*
- Example for model with two input features:

$$\mathcal{L}(\mathcal{D}; \mathbf{w}) = \sum_{n=1}^N \left[ q_1 q_2 + (1 - q_1) q_2 \exp(-y_n w_1 x_{n1}) \right. \\ \left. + (1 - q_2) q_1 \exp(-y_n w_2 x_{n2}) \right. \\ \left. + (1 - q_1)(1 - q_2) \exp(-y_n [w_1 x_{n1} + w_2 x_{n2}]) \right]$$

**Loss on first  
feature subset**



# Blankout: Ensemble interpretation

---

- MCF with blankout has an interesting interpretation as an *ensemble*
- Example for model with two input features:

$$\mathcal{L}(\mathcal{D}; \mathbf{w}) = \sum_{n=1}^N \left[ q_1 q_2 + (1 - q_1) q_2 \exp(-y_n w_1 x_{n1}) \right. \\ \left. + (1 - q_2) q_1 \exp(-y_n w_2 x_{n2}) \right. \\ \left. + (1 - q_1)(1 - q_2) \exp(-y_n [w_1 x_{n1} + w_2 x_{n2}]) \right]$$

**Loss on first  
feature subset**

**Loss on second  
feature subset**

# Blankout: Ensemble interpretation

---

- MCF with blankout has an interesting interpretation as an *ensemble*
- Example for model with two input features:

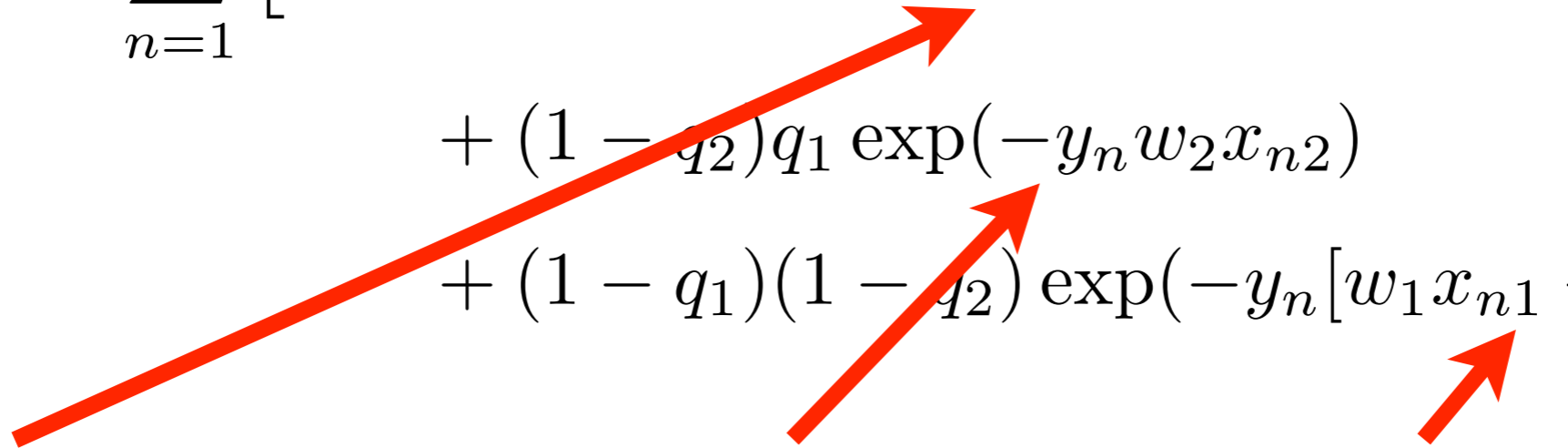
$$\mathcal{L}(\mathcal{D}; \mathbf{w}) = \sum_{n=1}^N \left[ q_1 q_2 + (1 - q_1) q_2 \exp(-y_n w_1 x_{n1}) \right. \\ \left. + (1 - q_2) q_1 \exp(-y_n w_2 x_{n2}) \right. \\ \left. + (1 - q_1)(1 - q_2) \exp(-y_n [w_1 x_{n1} + w_2 x_{n2}]) \right]$$

**Loss on first feature subset**      **Loss on second feature subset**      **Loss on full feature set**

# Blankout: Ensemble interpretation

---

- MCF with blankout has an interesting interpretation as an *ensemble*
- Example for model with two input features:

$$\mathcal{L}(\mathcal{D}; \mathbf{w}) = \sum_{n=1}^N \left[ q_1 q_2 + (1 - q_1) q_2 \exp(-y_n w_1 x_{n1}) \right. \\ \left. + (1 - q_2) q_1 \exp(-y_n w_2 x_{n2}) \right. \\ \left. + (1 - q_1)(1 - q_2) \exp(-y_n [w_1 x_{n1} + w_2 x_{n2}]) \right]$$


**Loss on first feature subset**      **Loss on second feature subset**      **Loss on full feature set**

- Note: MCF exponential loss is *convex* for all corrupting distributions



# Logistic loss

---

- Working out the MCF expectation (for independent corruption) gives:

$$\begin{aligned}\mathcal{L}(\mathcal{D}; \mathbf{w}) &= \sum_{n=1}^N \mathbb{E} \left[ \log \left( 1 + \exp \left( -y_n \mathbf{w}^T \tilde{\mathbf{x}}_n \right) \right) \right]_{p(\tilde{\mathbf{x}}_n | \mathbf{x}_n)} \\ &\leq \sum_{n=1}^N \log \left( 1 + \prod_{d=1}^D \mathbb{E} \left[ \exp \left( -y_n w_d \tilde{x}_{nd} \right) \right]_{p(\tilde{x}_{nd} | x_{nd})} \right)\end{aligned}$$

- The upper bound is obtained using *Jensen's inequality*\*

\* *Jensen's inequality*:  $\mathbb{E}[\phi(x)] \geq \phi(\mathbb{E}[x])$  for convex  $\phi(x)$

# Logistic loss

---

- Working out the MCF expectation (for independent corruption) gives:

$$\begin{aligned}\mathcal{L}(\mathcal{D}; \mathbf{w}) &= \sum_{n=1}^N \mathbb{E} \left[ \log \left( 1 + \exp \left( -y_n \mathbf{w}^T \tilde{\mathbf{x}}_n \right) \right) \right]_{p(\tilde{\mathbf{x}}_n | \mathbf{x}_n)} \\ &\leq \sum_{n=1}^N \log \left( 1 + \prod_{d=1}^D \mathbb{E} \left[ \exp \left( -y_n w_d \tilde{x}_{nd} \right) \right]_{p(\tilde{x}_{nd} | x_{nd})} \right)\end{aligned}$$

- The upper bound is obtained using *Jensen's inequality*\*
- Upper bound is *convex* iff the moment-generating function is *log-linear*

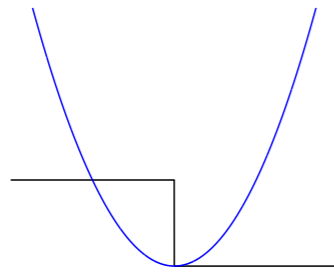
\* *Jensen's inequality*:  $\mathbb{E}[\phi(x)] \geq \phi(\mathbb{E}[x])$  for convex  $\phi(x)$

# Using MCF in practice

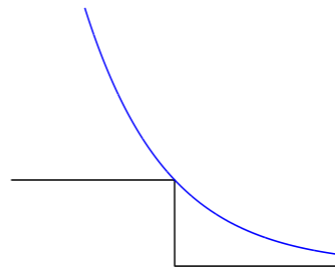
---

1)

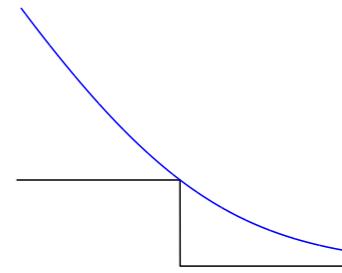
**Quadratic loss**



**Exponential loss**



**Logistic loss**

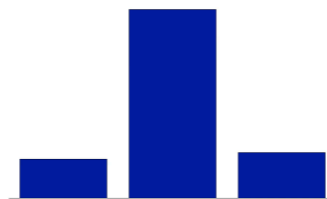


2)

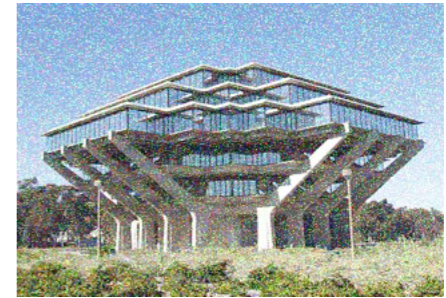
**Blankout noise**



**Poisson noise**



**Gaussian noise**



3)

**MCF Loss**

$$\mathcal{L}(\mathcal{D}; \mathbf{w})$$

# Experimental setup

---

- We performed three sets of experiments with MCF:
  - Document classification based on bag-of-word features
  - Image classification based on bag-of-visual-word features
  - *“Nightmare at test time”* scenario where features are unobserved at test time
- All our predictors use L2-regularization, with lambda set by cross-validation

# Experiment 1: Document classification

---

- We tested on three different document classification data sets
- All data sets have in the order of 20K features and 6K training examples

# Experiment 1: Document classification

---

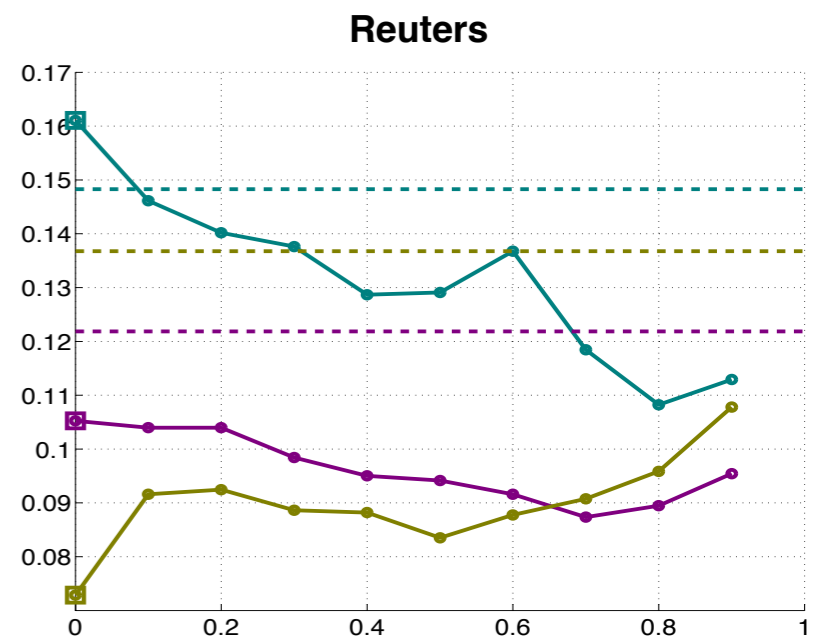
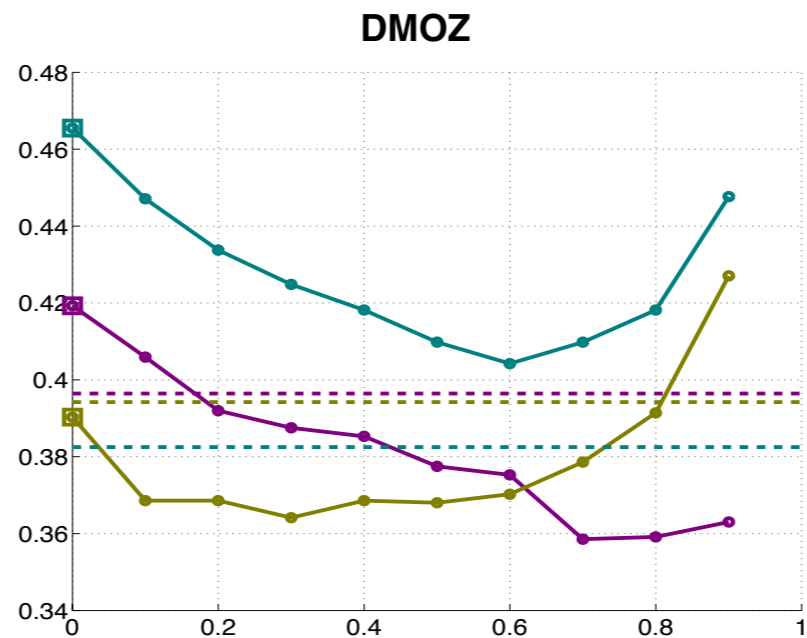
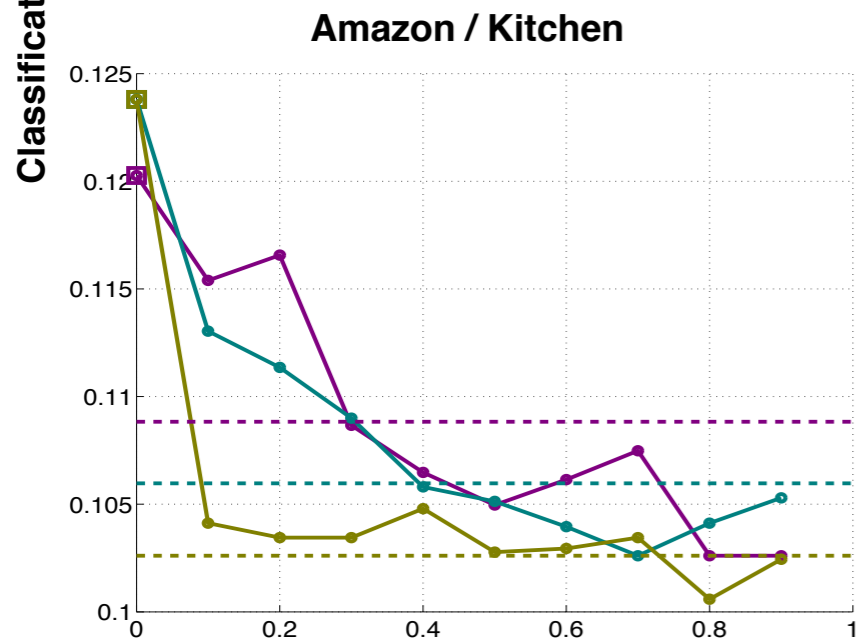
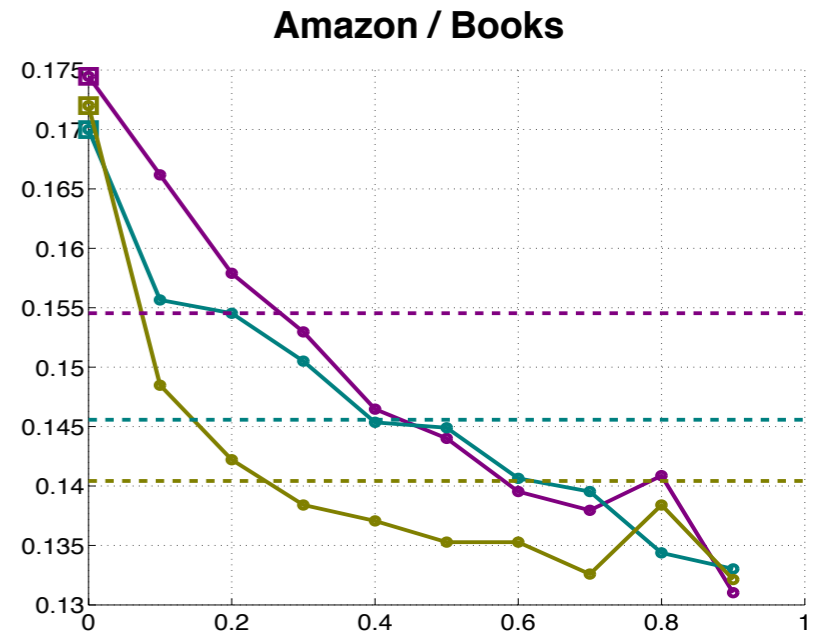
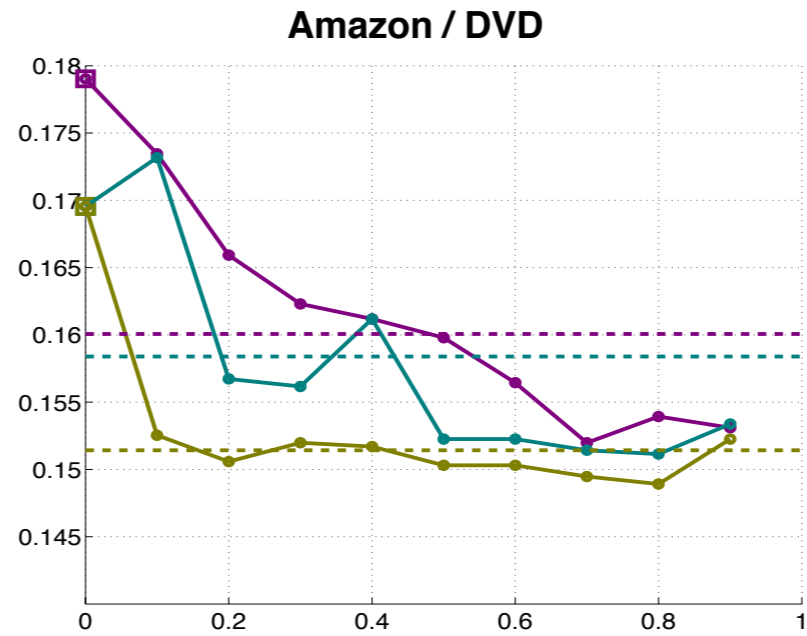
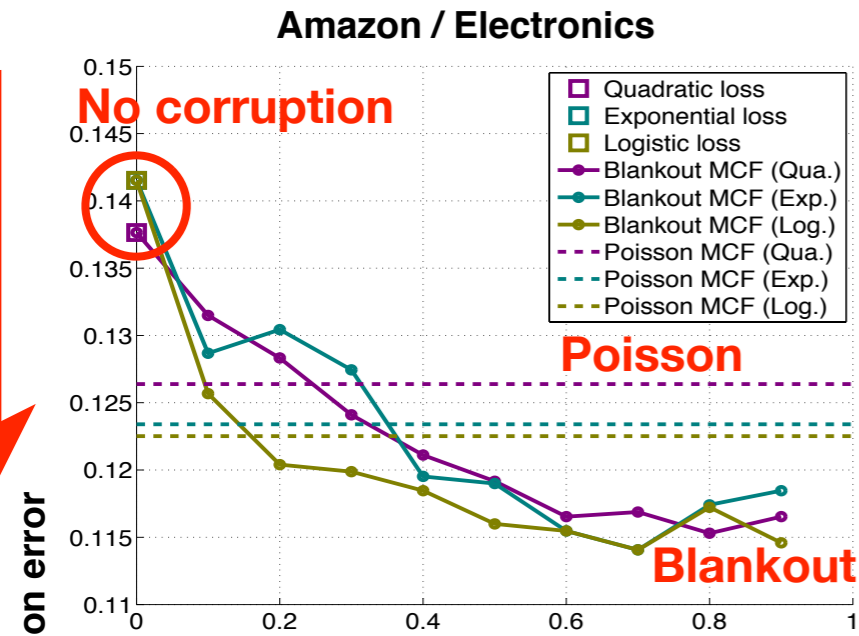
- We tested on three different document classification data sets
- All data sets have in the order of 20K features and 6K training examples
- We explore two different corrupting distributions:

- Blankout corruption: 
$$p(\tilde{x}_{nd} = 0) = q_d$$
$$p(\tilde{x}_{nd} = \frac{1}{1-q_d} x_{nd}) = 1 - q_d$$

- Poisson corruption: 
$$p(\tilde{x}_{nd}|x_{nd}) = \text{Pois}(\tilde{x}_{nd}|x_{nd})$$

# Experiment 1: Document classification

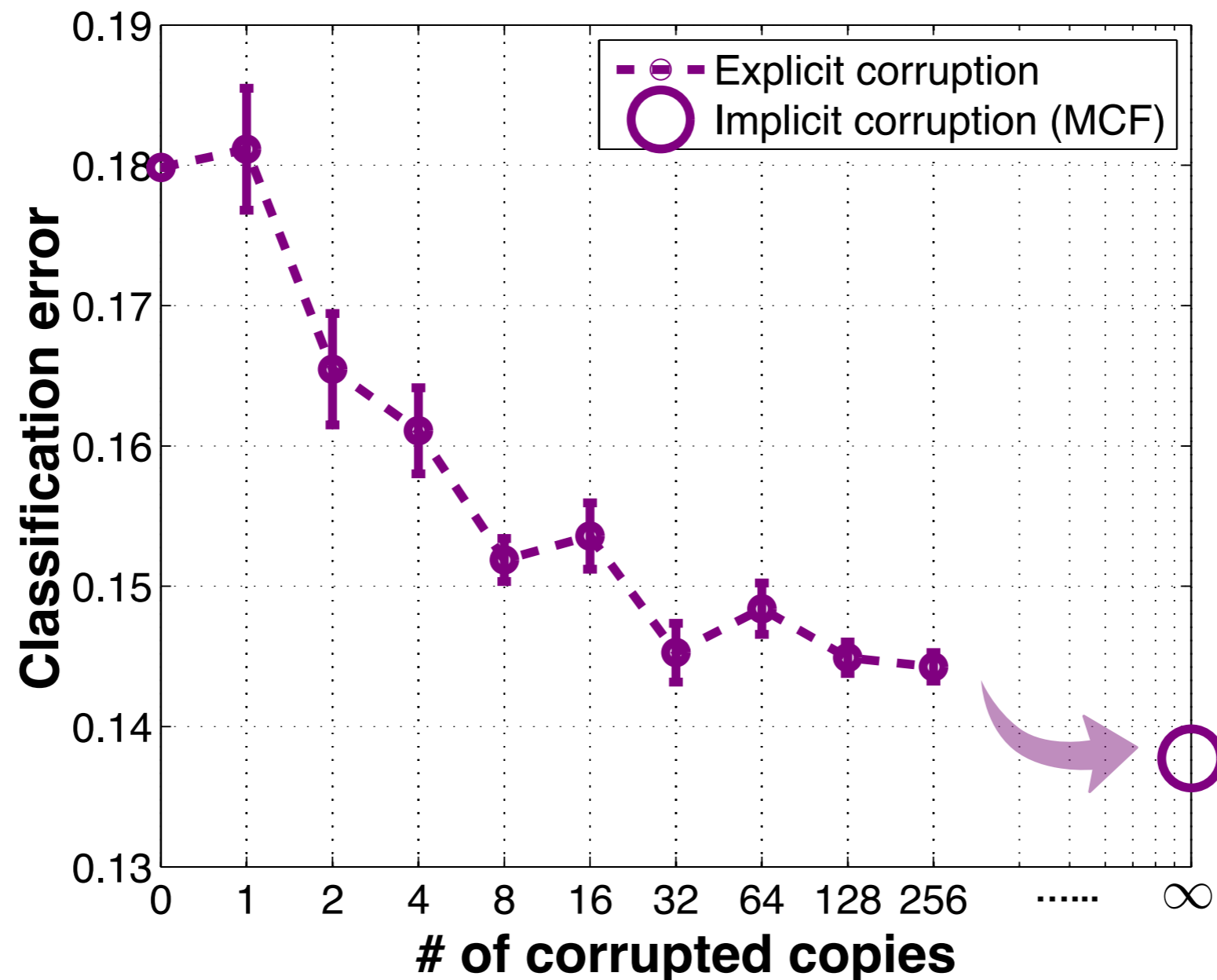
Better



More (blankout) corruption

# Experiment 1: Document classification

- Comparing *explicit* and *implicit* blankout corruption (Amazon Books; quadratic loss):

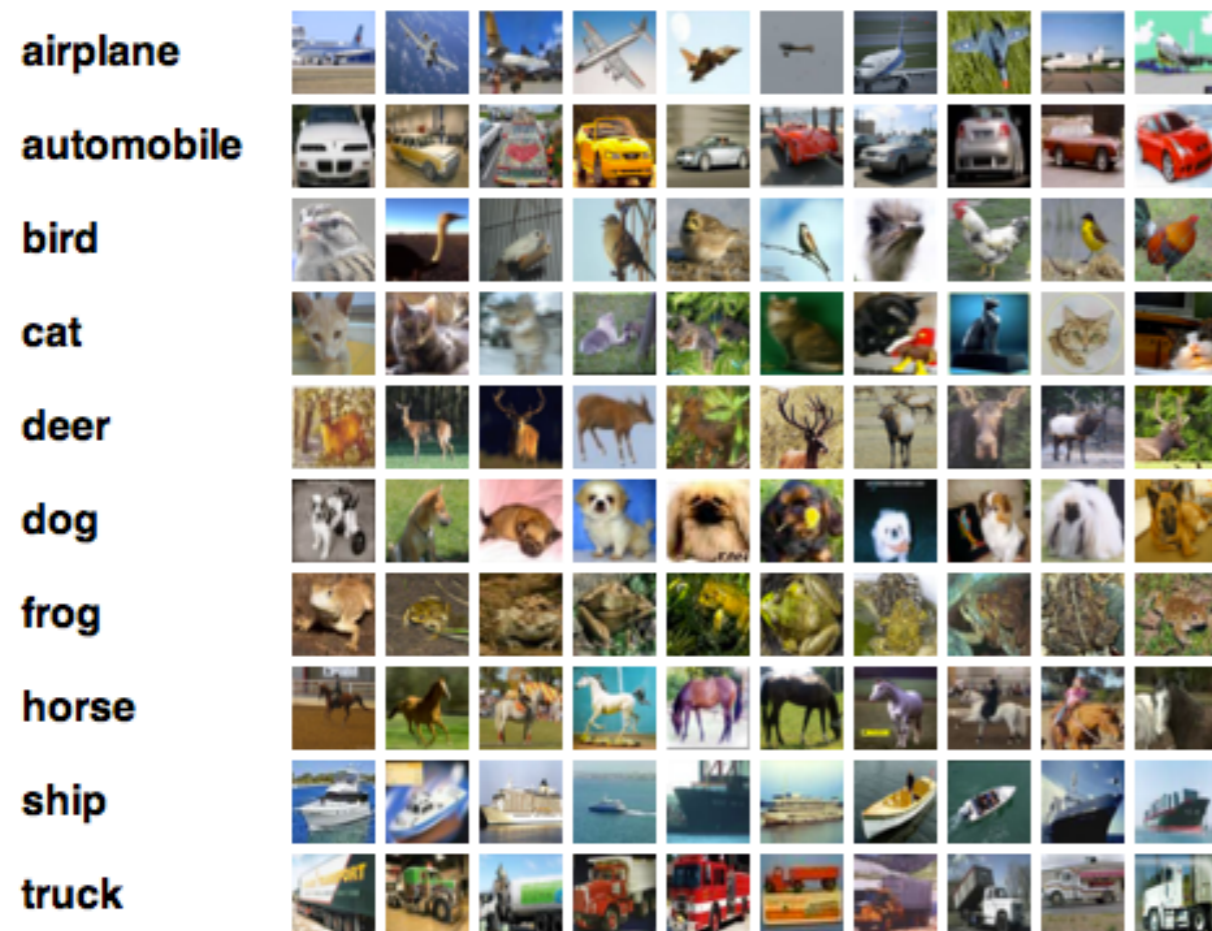




# Experiment 2: Image classification

---

- The CIFAR-10 data set contains 50K images of size 32x32 with 10 classes
- We use a standard\* bag-of-visual-words feature representation for the images



	<b>Quadr.</b>	<b>Expon.</b>	<b>Logist.</b>
<b>No MCF</b>	32.6%	39.7%	38.0%
<b>Poisson MCF</b>	29.1%	39.5%	30.0%
<b>Blankout MCF</b>	32.3%	37.9%	29.4%

\* We followed the approach by Coates et al. (2011) to extract features.

# Experiment 3: “Nightmare at test time”

---

- In some learning settings, features may be randomly unobserved at test time
- We experiment with this “nightmare at test time” scenario on MNIST digits:
  - Train regular and MCF-blankout classifiers on the original training set



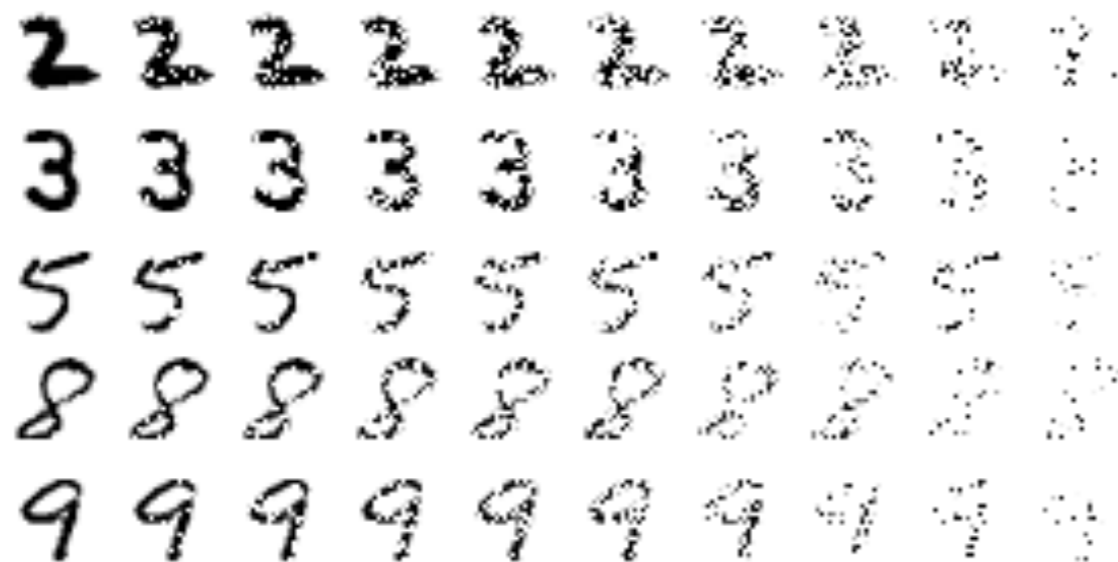
3	6	8	1	7	9	6	6	9	1
6	7	5	7	8	6	3	4	8	5
2	1	7	9	7	1	2	8	4	5
4	8	1	9	0	1	8	8	9	4
7	6	1	8	6	4	1	5	6	0
7	5	9	2	6	5	8	1	9	7
2	2	2	2	2	3	4	4	8	0
0	2	3	8	0	7	3	8	5	7
0	1	4	6	4	6	0	2	4	3
7	1	2	8	7	6	9	8	6	1

# Experiment 3: “Nightmare at test time”

---

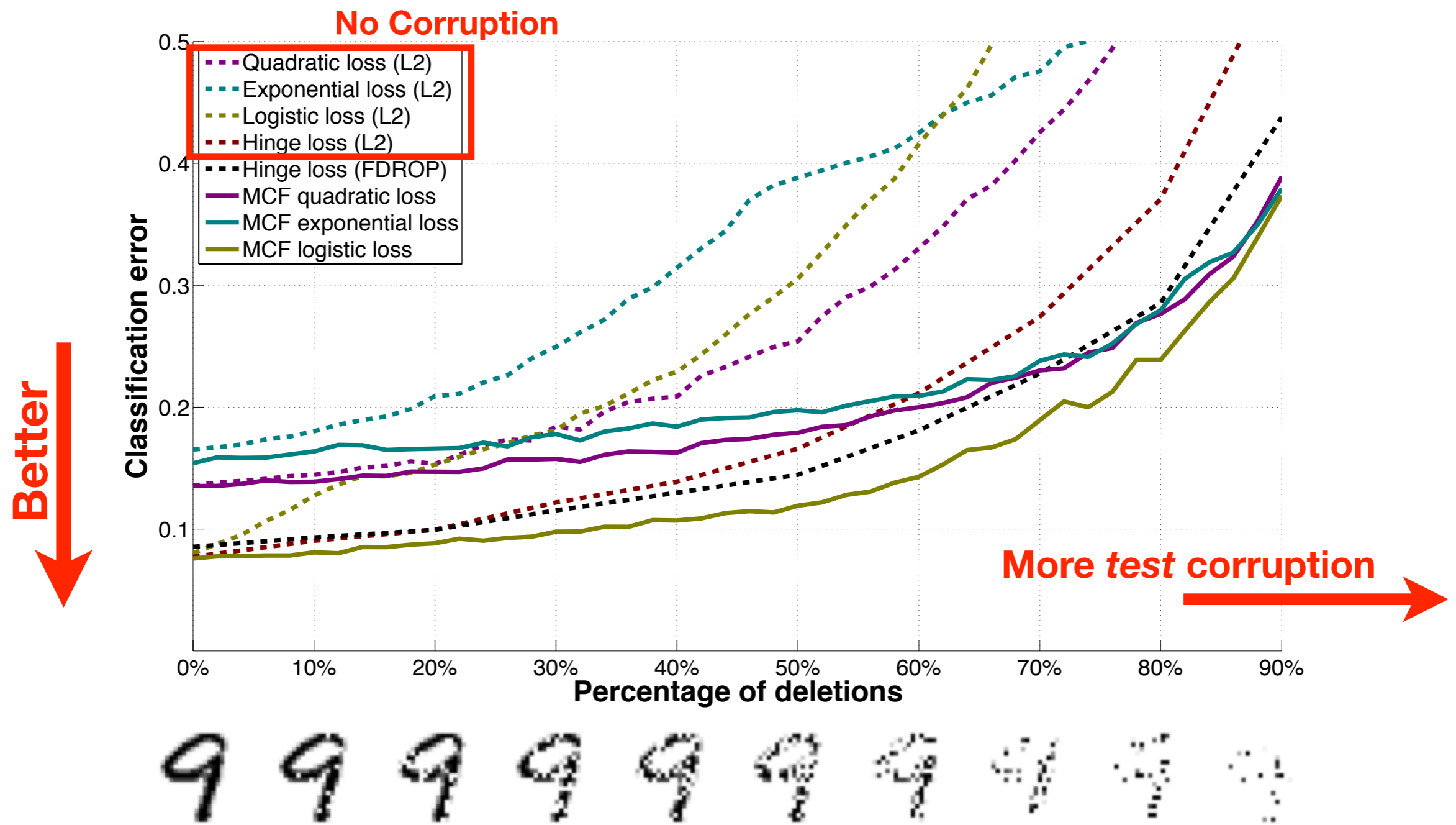
- In some learning settings, features may be randomly unobserved at test time
- We experiment with this “nightmare at test time” scenario on MNIST digits:

- Train regular and MCF-blankout classifiers on the original training set
- Randomly delete features from the test images, and measure classification error



# Experiment 3: “Nightmare at test time”

- Classification error on test images with randomly deleted features:



# Conclusions

---

- Adding *corrupted* examples to training data can *regularize* predictors

# Conclusions

---

- Adding *corrupted* examples to training data can *regularize* predictors
- For a range of models and corrupting distributions, MCF makes this *efficient*

# Conclusions

---

- Adding *corrupted* examples to training data can *regularize* predictors
- For a range of models and corrupting distributions, MCF makes this *efficient*
- MCF may lead to *improved results* in various learning settings:
  - In particular, in settings where you somewhat understand *how* data is generated
  - MCF may be very well suited for scenarios in which *domain shift* is present

# Thank you! Questions?

---

Thanks to:



Kilian Weinberger



Minmin Chen



Stephen Tyree