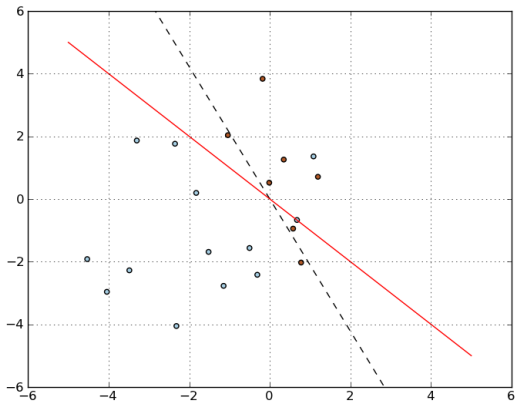# **Learning with Marginalized Corrupted Features**
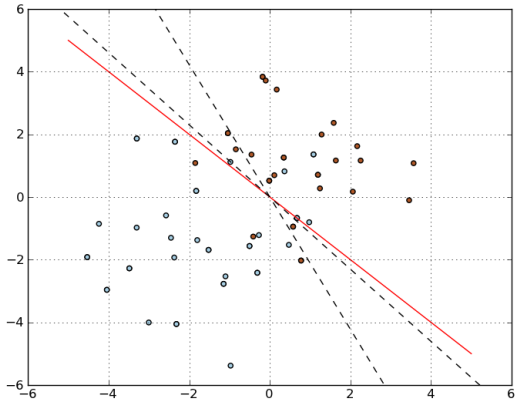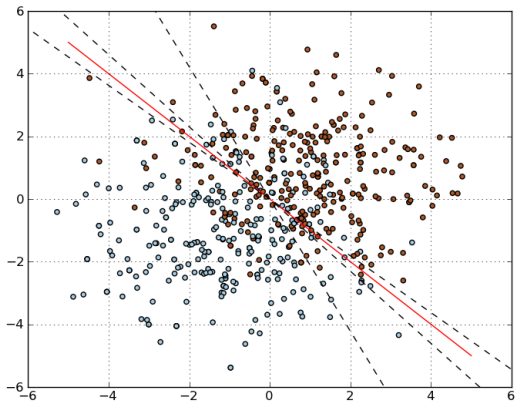
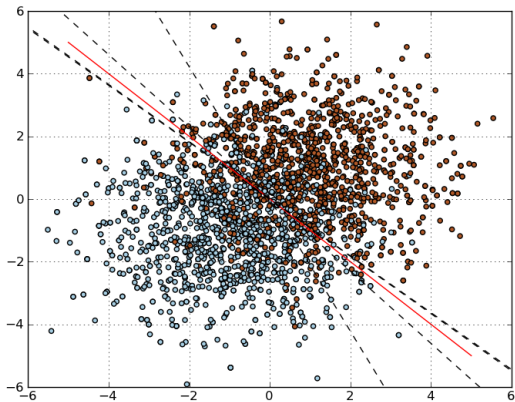L. van der Maaten, M. Chen, S. Tyree, K. Weinberger

ICML 2013

Jan Gasthaus

Tea talk
April 11, 2013

*Secret 4: lots of jittering, mirroring, and color perturbation of the original images generated on the fly to increase the size of the training set*

Yann LeCun on Google+ about Alex Krizhevsky's ImageNet results

**UCL**

- Old idea: create artificial additional training data by corrupting it with "noise"

- Old idea: create artificial additional training data by corrupting it with "noise"
- One easy way to incorporate domain knowledge (e.g. possible transformations)

- Old idea: create artificial additional training data by corrupting it with "noise"
- One easy way to incorporate domain knowledge (e.g. possible transformations)
- But: additional training data $\implies$ additional computation
- **Idea:** Corrupt with known ExpFam noise and integrate it out

- Explicit corruption: Take training set $D = \{(\mathbf{x}_n, y_n)\}_{n=1}^N$ and corrupt it $M$ times

$$\mathcal{L}(\tilde{D}, \Theta) = \sum_{n=1}^N \frac{1}{M} \sum_{m=1}^M L(\tilde{\mathbf{x}}_{nm}, y_n, \Theta)$$

with $\mathbf{x}_{nm} \sim p(\tilde{\mathbf{x}}_{nm}|\mathbf{x}_n)$.

- Implicit corruption: Minimize the expected value of the loss under $p(\tilde{\mathbf{x}}_n|\mathbf{x}_n)$:

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

i.e. replace the empirical average with the expectation.

**UCL**

- This is *so* obvious that it must have been done before . . .

- This is *so* obvious that it must have been done before . . .
  - *Vicinal Risk Minimization*, Chapelle, Weston, Bottou, & Vapnik, NIPS 2000

$$R_{vic}(f) = \int \ell(f(\mathbf{x}), y) \, dP_{est}(\mathbf{x}, y) = \frac{1}{n} \sum_{i=1}^{n} \int \ell(f(\mathbf{x}), y_i) dP_{\mathbf{x}_i}(\mathbf{x})$$

- This is *so* obvious that it must have been done before . . .
  - *Vicinal Risk Minimization*, Chapelle, Weston, Bottou, & Vapnik, NIPS 2000

$$R_{vic}(f) = \int \ell(f(\mathbf{x}), y) \, dP_{est}(\mathbf{x}, y) = \frac{1}{n} \sum_{i=1}^{n} \int \ell(f(\mathbf{x}), y_i) dP_{\mathbf{x}_i}(\mathbf{x})$$

- Explicitly only consider the case of Gaussian noise distributions

**Quadratic loss.** Assuming[2] a label variable $y \in \{-1, +1\}$, the expected value of the quadratic loss under corrupting distribution $p(\tilde{\mathbf{x}}|\mathbf{x})$ is given by:

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

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

$$- 2\left(\sum_{n=1}^{N} y_n \mathbb{E}[\tilde{\mathbf{x}}_n]\right)^{\mathrm{T}} \mathbf{w} + N, \quad (3)$$

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

# Quadratic Loss

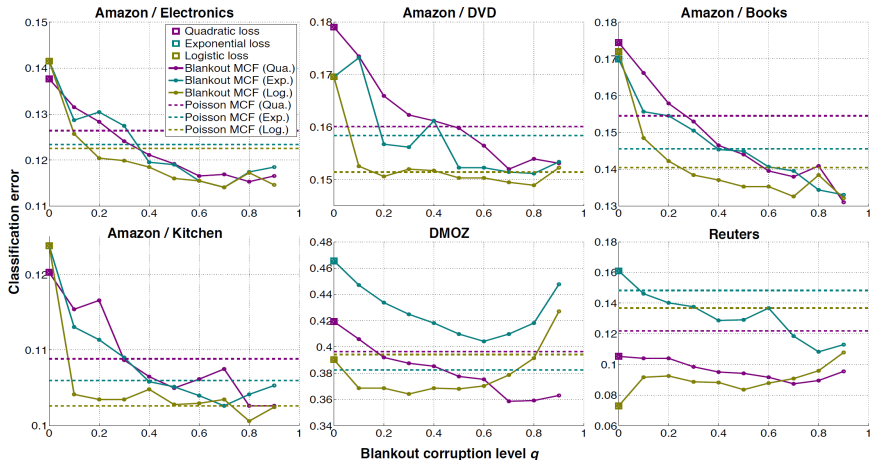| Distribution | PDF | $\mathbb{E}[\tilde{x}_{nd}]_{p(\tilde{x}_{nd}|x_{nd})}$ | $V[\tilde{x}_{nd}]_{p(\tilde{x}_{nd}|x_{nd})}$ |
|---|---|---|---|
| Blankout noise | $p(\tilde{x}_{nd} = 0) = q_d$ <br> $p(\tilde{x}_{nd} = \frac{1}{1-q_d} x_{nd}) = 1 - q_d$ | $x_{nd}$ | $\frac{q_d}{1-q_d} x_{nd}^2$ |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

**Exponential loss.** The expected value of the exponential loss under corruption model $p(\tilde{\mathbf{x}}|\mathbf{x})$ is:
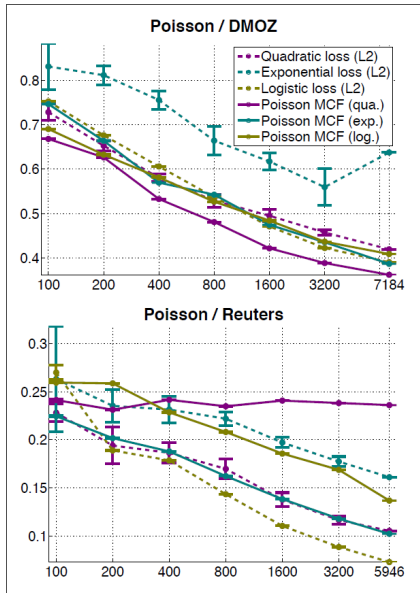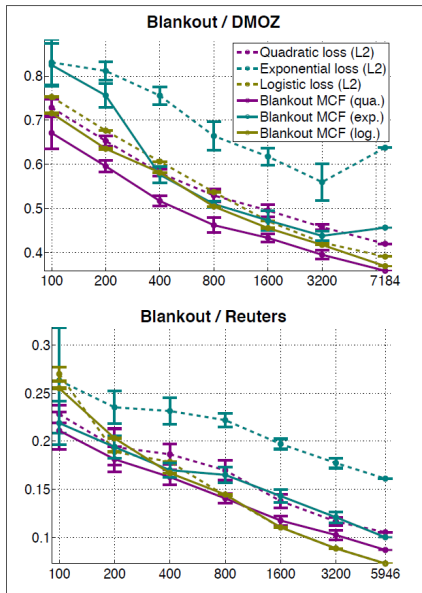
$$\mathcal{L}(\mathcal{D}; \mathbf{w}) = \sum_{n=1}^{N} \mathbb{E}\left[e^{-y_n \mathbf{w}^{\mathrm{T}} \tilde{\mathbf{x}}_n}\right]_{p(\tilde{\mathbf{x}}_n | \mathbf{x}_n)}$$

$$= \sum_{n=1}^{N} \prod_{d=1}^{D} \mathbb{E}\left[e^{-y_n w_d \tilde{x}_{nd}}\right]_{p(\tilde{x}_{nd} | x_{nd})}, \quad (4)$$

**Logistic loss.** In the case of the logistic loss, the solution to (2) cannot be computed in closed form. Instead, we derive an upper bound, which can be minimized as a surrogate loss:

$$\mathcal{L}(\mathcal{D}; \mathbf{w}) = \sum_{n=1}^{N} \mathbb{E}\left[\log\left(1 + e^{-y_n \mathbf{w}^\mathrm{T} \tilde{\mathbf{x}}_n}\right)\right]_{p(\tilde{\mathbf{x}}_n | \mathbf{x}_n)}$$

$$\leq \sum_{n=1}^{N} \log\left(1 + \prod_{d=1}^{D} \mathbb{E}\left[e^{-y_n w_d \tilde{x}_{nd}}\right]_{p(\tilde{x}_{nd} | x_{nd})}\right). \quad (5)$$

| **Distribution** | $\mathbb{E}[\exp(-y_n w_d \tilde{x}_{nd})]_{p(\tilde{x}_{nd}\mid x_{nd})}$ |
|---|---|
| Blankout noise | $q_d + (1 - q_d)\exp(-y_n w_d \frac{1}{1-q_d} x_{nd})$ |
| Gaussian noise | $\exp(-y_n w_d x_{nd} + \frac{1}{2}\sigma^2 y_n^2 w_d^2)$ |
| Laplace noise | $(1 - \lambda^2 y_n^2 w_d^2)^{-1}\exp(-y_n w_d x_{nd})$ |
| Poisson noise | $\exp(x_{nd}(\exp(-y_n w_d) - 1))$ |

# Results

# Results

**Blankout / DMOZ**

Quadratic loss (L2)
Exponential loss (L2)
Logistic loss (L2)
Blankout MCF (qua.)
Blankout MCF (exp.)
Blankout MCF (log.)

**Poisson / DMOZ**

Quadratic loss (L2)
Exponential loss (L2)
Logistic loss (L2)
Poisson MCF (qua.)
Poisson MCF (exp.)
Poisson MCF (log.)

**Blankout / Reuters**

**Poisson / Reuters**

| | Quadr. | Expon. | Logist. |
|---|---|---|---|
| **No MCF** | 32.6% | 39.7% | 32.5% |
| **Poisson MCF** | 29.1% | 39.5% | 30.0% |
| **Blankout MCF** | 32.3% | 37.9% | 29.4% |

*Table 3.* Classification errors obtained on the CIFAR-10 data set with MCF classifiers trained on simple spatial-pyramid bag-of-visual-words features.