

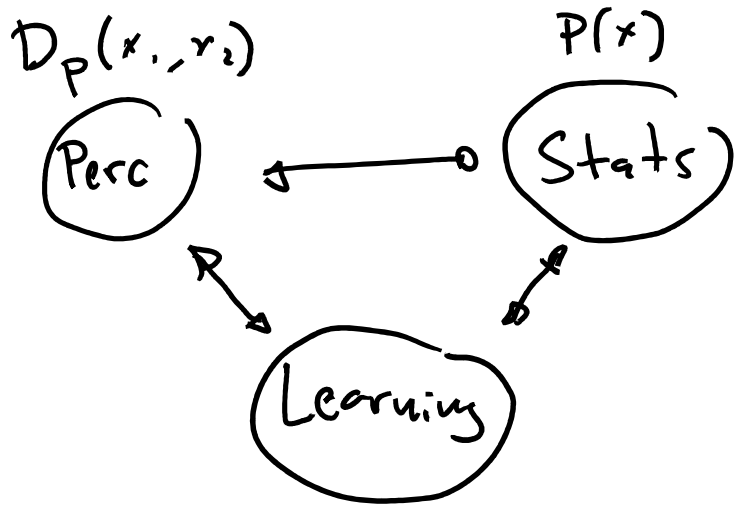
SOME OBSERVATIONS ON THE RELATION BETWEEN
PERCEPTUAL DISTANCE AND STATISTICAL LEARNING
IN AUTOENCODERS

A. Hejblum, V. Laparra, R. Santos, J. Ballé & J. Malo

<https://arxiv.org/abs/2106.04427>

7th july 2021

ON THE RELATION BETWEEN STAT. LEARNING AND PERCEPTUAL DISTANCE



Autoencoders



$$f = d \circ e$$

* PERCEPTION: $x \xrightarrow{S} y$

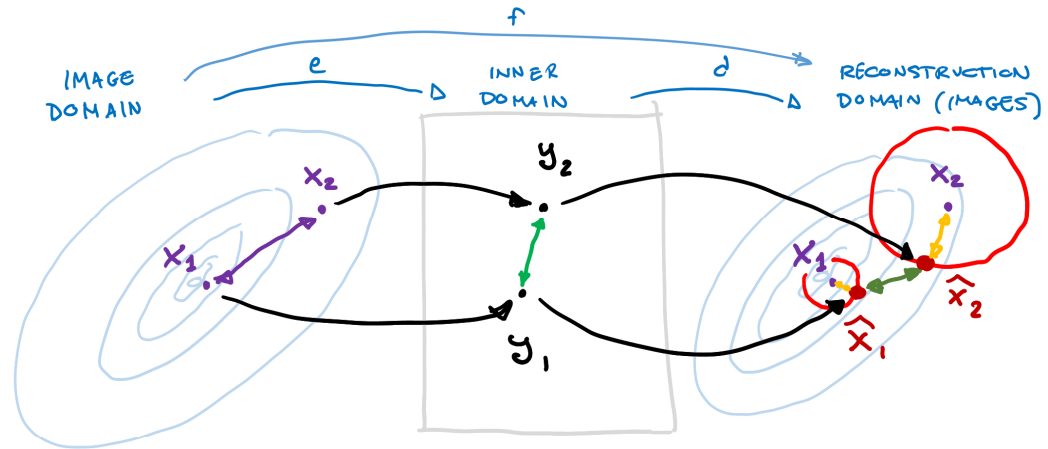
$$D_p(x, x_2) = \|y_1 - y_2\|_2 = \|S(x_1) - S(x_2)\|_2$$

* LEARNING

- Optimization $f_\theta(x)$

$$\mathcal{L}(\theta) = \mathbb{E}_x [D(x, f_\theta(x))] = \int p(x) D(x, f(x)) dx$$

- Induced distances?



EUCLIDEAN DISTANCE $D_e(x_1, x_2) = \|x_1 - x_2\|_2 = \text{RMSE}(x_1, x_2)$

INNER DISTANCE $D_i(x_1, x_2) = \|y_1 - y_2\|_2 = \|e(x_1) - e(x_2)\|_2$

RECONSTRUCT. DISTANCE $D_r(x_1, x_2) = \|\hat{x}_1 - \hat{x}_2\|_2 = \|f(x_1) - f(x_2)\|_2$

SELF RECONSTRUCT. DISTANCE $D_s(x, \hat{x}) = \|x - \hat{x}\|_2 = \|x - f(x)\|_2$

Proof of conjecture 1 (sensitivity of D_p) assuming S equalizes $p(x)$

Obs. 1 Sensitivity $D_p \sim p(x)$ $D_p = (\Delta S^T \Delta S)^{\frac{1}{2}}$

(A) $D_p(x_1, x_2) = |S(x_1) - S(x_2)|_2 = |S(x_1) - S(x_1 + \Delta x)|_2 = |S(x_1) - S(x_1) - \frac{\partial S}{\partial x} \Delta x|_2$

$D_p(x_1, x_2) = \left| \frac{\partial S}{\partial x} \Delta x \right|_2 = \left(\Delta x^T \left[\frac{\partial S}{\partial x} \right]^T \left[\frac{\partial S}{\partial x} \right] \Delta x \right)^{\frac{1}{2}} = \frac{\partial S}{\partial x} \Delta x \Rightarrow \frac{\partial S}{\partial x} = \frac{D_p}{\Delta x}$

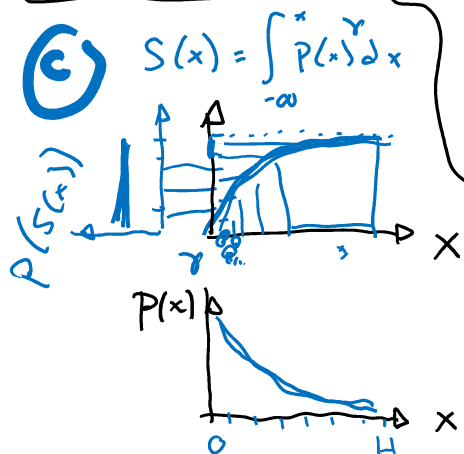
$\Delta S^T \cdot \Delta S$ ΔS
 Δx^T Δx Δx Δx
 D_p^2 Δ RMSE *

(B) Sensitivity $\frac{\partial D_p}{\partial x}$

$\frac{\partial D_p}{\partial x} = \frac{\partial}{\partial x} (\Delta S^T \Delta S)^{\frac{1}{2}} = \frac{\partial (\Delta S^T \Delta S)^{\frac{1}{2}}}{\partial (\Delta S^T \Delta S)} \cdot \frac{\partial (\Delta S^T \Delta S)}{\partial \Delta S} \cdot \frac{\partial \Delta S}{\partial x}$

$= \frac{1}{2} (\Delta S^T \Delta S)^{-\frac{1}{2}} \cdot \Delta S^T \cdot \frac{\partial S}{\partial x} = \frac{1}{D_p} \cdot \Delta S^T \cdot \frac{\partial S}{\partial x}$

$\frac{\partial D_p}{\partial x} \stackrel{\Delta}{=} \frac{1}{D_p} \Delta S \frac{\partial S}{\partial x} \Rightarrow \frac{\partial D_p}{\partial x} = \frac{\partial S}{\partial x} **$



$S(x) = \int_{-\infty}^x p(x) dx \Rightarrow \frac{\partial S}{\partial x} = p(x) ***$

Laughlin 81, McLeod 01, Malo 06, Laparra 12, Laparra 15

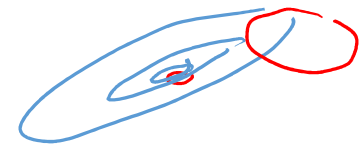
$\frac{\partial D_p}{\partial x} = \frac{D_p}{RMSE} = p(x)$
 SENSITIVITY OF D_p

Obs. 2 Self reconst. distance $D_S \propto \frac{1}{P(x)}$

$$\mathcal{L}(\theta) = \int p(x) \underbrace{D(x, f_\theta(x))}_{\substack{\uparrow \\ D \equiv \text{RMSE} \\ \text{MSE}}} dx = \int p(x) \underbrace{|x - f(x)|_2}_{D_S} dx \rightarrow$$

$$D_S \propto \frac{1}{P(x)^2}$$

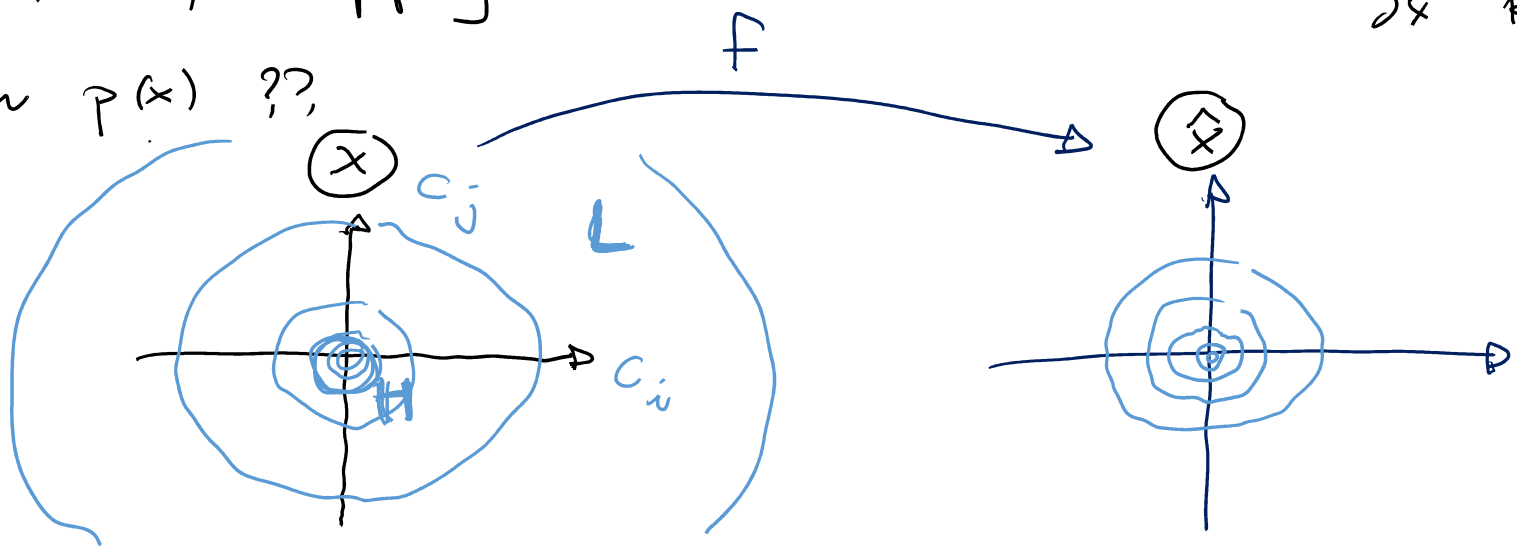
$$D_S = |x - \hat{x}|_2$$

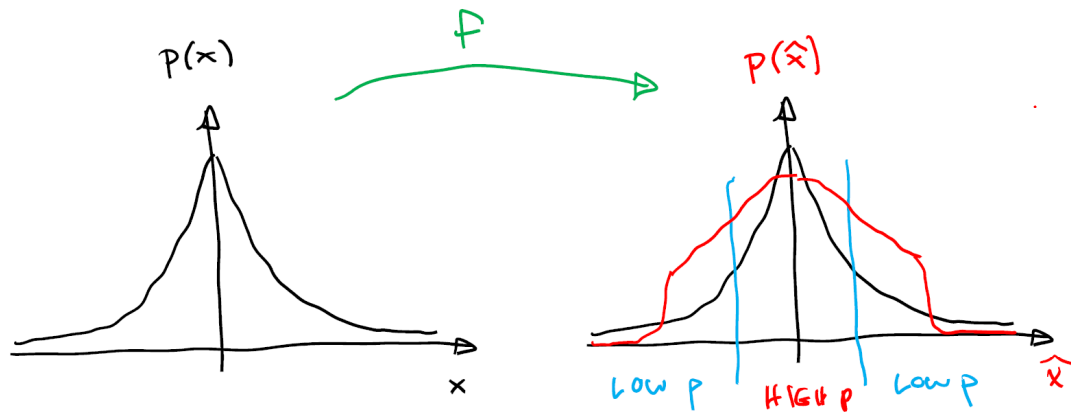


Obs. 3 sensit. of $D_r \propto P(x)$

Why? \bullet A, B apply to f (as well as S') $\Rightarrow \frac{\partial D_r}{\partial x} = \frac{D_r}{\text{RMSE}} = \frac{\partial f}{\partial x}$

$\bullet \frac{\partial f}{\partial x} \sim P(x) ??$





f is "kind of" equalizing...
 Not totally, but a little bit!!

$$\Rightarrow \frac{\partial f}{\partial x} \sim p(x)^\gamma \Rightarrow \frac{\partial D_r}{\partial x} \propto p(x)$$

Observat. 3
 Sensitivity of D_r

$$\frac{\partial D_r}{\partial x} = \frac{|f(x_1) - f(x_2)|_2}{|x_1 - x_2|_2} \propto p(x)^\gamma$$

Observation 4: Sensitivity of D_i

$$\frac{\partial D_i}{\partial x} = \frac{|e(x_1) - e(x_2)|_2}{|x_1 - x_2|_2} \propto p(x)^\gamma$$

Conjectures / observations

① Sensitivity of D_p _____

$$\frac{\partial D_p}{\partial x} \equiv \frac{|S(x_1) - S(x_2)|_2}{|x_1 - x_2|_2} = P(x)^r$$

② Self reconstruct. distance D_s _____

$$D_s(x, \hat{x}) \propto \frac{1}{P(x)}$$

③ Sensitivity of D_r _____

$$\frac{\partial D_r}{\partial x} \equiv \frac{|f(x_1) - f(x_2)|_2}{|x_1 - x_2|_2} \propto P(x)$$

④ Sensitivity of D_i _____

$$\frac{\partial D_i}{\partial x} \equiv \frac{|e(x_1) - e(x_2)|_2}{|x_1 - x_2|_2} \propto P(x)$$

Conjectures / observations

① Sensitivity of D_p $\frac{\partial D_p}{\partial x} \equiv \frac{|S(x_1) - S(x_2)|_2}{|x_1 - x_2|_2} = p(x)^r$

Equalization Hypothesis

Laughlin 81, McLeod 01, Malo 06, Laparra 12, Laparra 15

Malo 20 $I(x, S(x)) = \sum_i \underline{h(s(x)_i)} - \underline{T(s(x))} - \underline{h(u)}$
 Equalization is good. Redundancy & Noise are bad!

② Self reconstruct. distance D_s $D_s(x, \hat{x}) \propto \frac{1}{p(x)}$

③ Sensitivity of D_r $\frac{\partial D_r}{\partial x} \equiv \frac{|f(x_1) - f(x_2)|_2}{|x_1 - x_2|_2} \propto p(x)$

④ Sensitivity of D_i $\frac{\partial D_i}{\partial x} \equiv \frac{|e(x_1) - e(x_2)|_1}{|x_1 - x_2|_2} \propto p(x)$

Conjectures / observations

① Sensitivity of D_p $\frac{\partial D_p}{\partial x} \equiv \frac{|S(x_1) - S(x_2)|_2}{|x_1 - x_2|_2} = p(x)^r$

Equalization Hypothesis

Laughlin 81, McLeod 01, Malo 06, Laparra 12, Laparra 15

Malou 20 $I(x, S(x)) = \sum_i \underbrace{h(s(x)_i)} - \underbrace{T(s(x))} - \underbrace{h(u)}$
 Equalization is good. Redundancy & Noise are bad!

② Self reconstruct. distance D_s $D_s(x, \hat{x}) \propto \frac{1}{p(x)}$

$\chi(\theta) = \int p(x) D(x, f_\theta(x)) dx$

③ Sensitivity of D_r $\frac{\partial D_r}{\partial x} \equiv \frac{|f(x_1) - f(x_2)|_2}{|x_1 - x_2|_2} \propto p(x)$

④ Sensitivity of D_i $\frac{\partial D_i}{\partial x} \equiv \frac{|e(x_1) - e(x_2)|_1}{|x_1 - x_2|_2} \propto p(x)$

Conjectures / observations

① Sensitivity of D_p

$$\frac{\partial D_p}{\partial x} \equiv \frac{|S(x_1) - S(x_2)|_2}{|x_1 - x_2|_2} = p(x)^r$$

Equalization Hypothesis

Laughlin 81, McLeod 01, Malo 06, Laparra 12, Laparra 15

Malu 20 $I(x, S(x)) = \sum_i \underline{h(s(x)_i)} - \underline{T(s(x))} - \underline{h(u)}$
 Equalization is good. Redundancy & Noise are bad!

② Self reconstruct. distance D_S

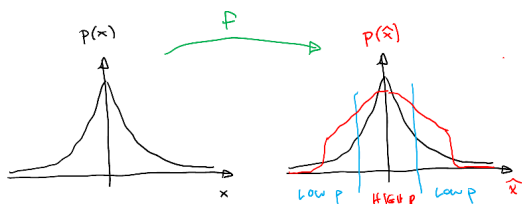
$$D_S(x, \hat{x}) \propto \frac{1}{p(x)}$$

$$\chi(\theta) = \int p(x) D(x, f_\theta(x)) dx$$

③ Sensitivity of D_r

$$\frac{\partial D_r}{\partial x} \equiv \frac{|f(x_1) - f(x_2)|_2}{|x_1 - x_2|_2} \propto p(x)$$

f kind of equalizer



④ Sensitivity of D_i

$$\frac{\partial D_i}{\partial x} \equiv \frac{|e(x_1) - e(x_2)|_1}{|x_1 - x_2|_2} \propto p(x)$$

Conjectures / observations

① Sensitivity of D_p

$$\frac{\partial D_p}{\partial x} \equiv \frac{|S(x_1) - S(x_2)|_2}{|x_1 - x_2|_2} = p(x)^r$$

Equalization Hypothesis

Laughlin 81, McLeod 01, Malo 06, Laparra 12, Laparra 15

$$\left[\begin{array}{l} \text{Malo 20} \\ \text{Equalization is good, Redundancy \& Noise are bad!} \end{array} \right. \quad I(x, S(x)) = \sum_i \underbrace{h(S(x)_i)} - \underbrace{T(S(x))} - \underbrace{h(u)}$$

② Self reconstruct. distance D_S

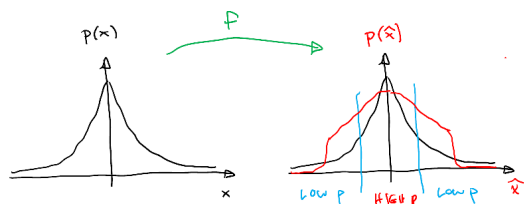
$$D_S(x, \hat{x}) \propto \frac{1}{p(x)}$$

$$\chi(\theta) = \int p(x) D(x, f_\theta(x)) dx$$

③ Sensitivity at D_r

$$\frac{\partial D_r}{\partial x} \equiv \frac{|f(x_1) - f(x_2)|_2}{|x_1 - x_2|_2} \propto p(x)$$

f kind of equalizer



④ Sensitivity of D_i

$$\frac{\partial D_i}{\partial x} \equiv \frac{|e(x_1) - e(x_2)|_1}{|x_1 - x_2|_2} \propto p(x)$$

The inner domain maybe is equalizing

Conjectures / observations

$D_p \equiv$ Div. Norm. (NLPD) / MS-SSIM
 $p(x) \equiv$ PIXEL CNN + <https://arxiv.org/abs/1701.05517>

(1) Sensitivity of D_p

$$\frac{\partial D_p}{\partial x} \equiv \frac{|S(x_1) - S(x_2)|_2}{|x_1 - x_2|_2} = p(x)^r$$

Equalization Hypothesis

Laughlin 81, McLeod 01, Malo 06, Laparra 12, Laparra 15

Malu 20 $I(x, S(x)) = \sum_i \underbrace{h(S(x)_i)} - \underbrace{T(S(x))} - \underbrace{h(u)}$
 Equalization is good. Redundancy & Noise are bad!

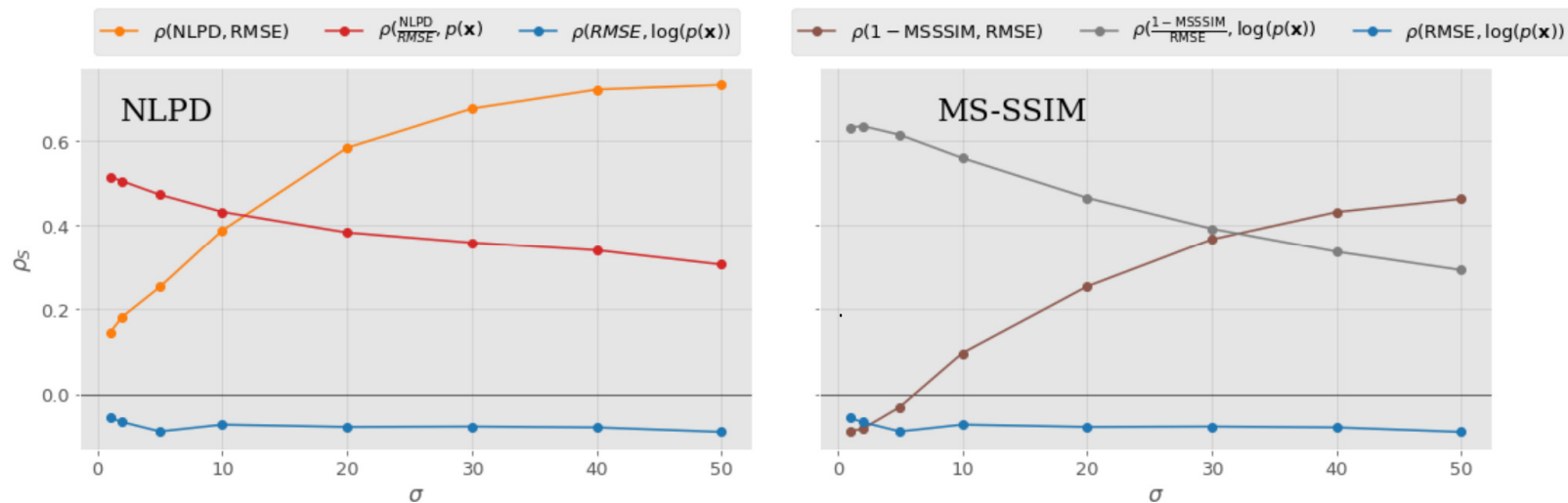


Figure 2: **Illustration of Observation 1.** Spearman correlations ρ_S between the sensitivity of the perceptual distances NLPD and MS-SSIM and $\log(p(x))$ (in red/gray). Distances are computed between x , and a distorted version with additive Gaussian noise, $x + \Delta x$, with deviation σ . Correlation of RMSE with perceptual distortions (in orange/brown) and of RMSE with $\log(p(x))$ (in blue) are included for comparison. MS-SSIM is a similarity index, so 1-(MS-SSIM) is a distortion measure.

Conjectures / observations $\mathcal{L} = E_x [H(\hat{y}) + \lambda D_e(x, \hat{x})]$ <https://arxiv.org/pdf/2007.03034.pdf>

- (2) Self reconstruct. distance D_s ————— $D_s(x, \hat{x}) \propto \frac{1}{p(x)}$ Empirical Risk Minimize
- (3) Sensitivity at D_r ————— $\frac{\partial D_r}{\partial x} \equiv \frac{|f(x_1) - f(x_2)|_2}{|x_1 - x_2|_2} \propto p(x)$ Equalized?
- (4) Sensitivity of D_i ————— $\frac{\partial D_i}{\partial x} \equiv \frac{|e(x_1) - e(x_2)|_1}{|x_1 - x_2|_2} \propto p(x)$ Equalized?

https://github.com/tensorflow/compression/blob/master/models/toy_sources/toy_sources.ipynb

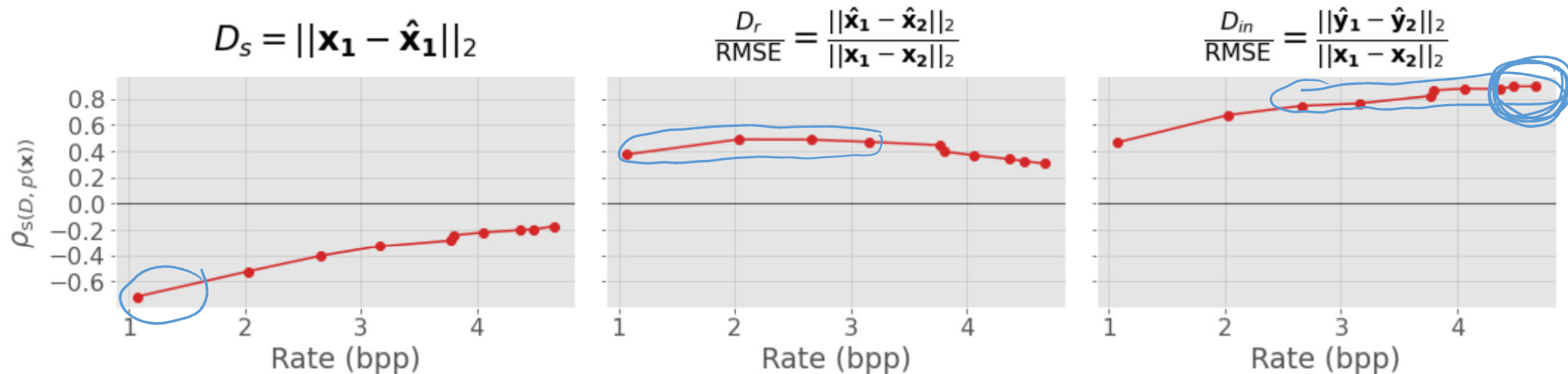


Figure 3: **Illustrating Observations 2, 3, and 4.** Spearman correlations $\rho_S(D, p(x_1))$ between different distances and the probability of point x_1 are shown. Each point corresponds to the correlation for one autoencoder trained for a particular Rate regime.

Conjectures / Observations

(5) The induced distances D_r and D_i are related to D_p

$$\begin{array}{l} \textcircled{1} \text{ Sensitivity of } D_p \quad \frac{\partial D_p}{\partial x} = \frac{|S(x_1) - S(x_2)|_2}{|x_1 - x_2|_2} = P(x)^T \\ \textcircled{3} \text{ Sensitivity of } D_r \quad \frac{\partial D_r}{\partial x} = \frac{|f(x_1) - f(x_2)|_2}{|x_1 - x_2|_2} \propto P(x) \\ \textcircled{4} \text{ Sensitivity of } D_i \quad \frac{\partial D_i}{\partial x} = \frac{|e(x_1) - e(x_2)|_1}{|x_1 - x_2|_2} \propto P(x) \end{array} \quad \left. \vphantom{\begin{array}{l} \textcircled{1} \\ \textcircled{3} \\ \textcircled{4} \end{array}} \right\} \Rightarrow \begin{array}{l} D_i \propto D_p \text{ (MOS)} \\ D_r \propto D_p \text{ (MOS)} \end{array}$$

(6) As D_p is related to $P(x)$, \Rightarrow using $\mathcal{L}(\theta) = \int P(x) D_p(x, f(x)) dx$ implies considering $P(x)$ twice

(6.A) GOOD When your data is poor using D_p regularizes

(6.B) BAD When training with D_p performance may be poor in $P(x)$ \downarrow

Conjectures / Observations

(5) The induced distances D_r and D_i are related to D_p

$$\begin{aligned}
 \textcircled{1} \text{ Sensitivity of } D_p & \quad \frac{\partial D_p}{\partial x} = \frac{|S(x_1) - S(x_2)|_2}{|x_1 - x_2|_2} = p(x)^Y \\
 \textcircled{3} \text{ Sensitivity of } D_r & \quad \frac{\partial D_r}{\partial x} = \frac{|f(x_1) - f(x_2)|_2}{|x_1 - x_2|_2} \propto p(x) \\
 \textcircled{4} \text{ Sensitivity of } D_i & \quad \frac{\partial D_i}{\partial x} = \frac{|e(x_1) - e(x_2)|_1}{|x_1 - x_2|_2} \propto p(x)
 \end{aligned}
 \quad \Rightarrow \quad
 \begin{aligned}
 D_i & \propto D_p \text{ (MOS)} \\
 D_r & \propto D_p \text{ (MOS)}
 \end{aligned}$$

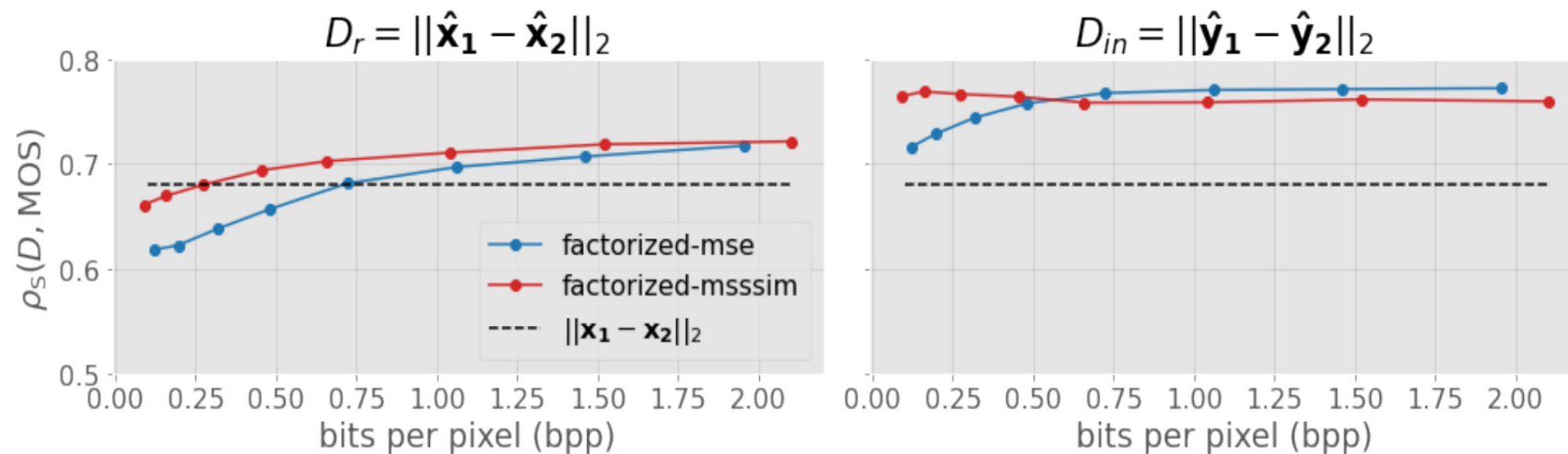
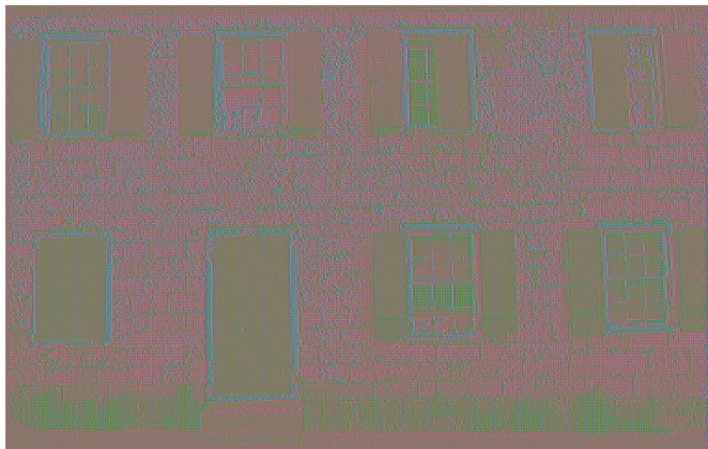


Figure 4: **Checking Observation 5.** Spearman correlations ρ_S between induced distances (D_r or D_{in}) and mean opinion score (MOS) for images from TID 2013 dataset [33]. Pretrained compressive autoencoders at different bitrates were used. *factorized-mse* denotes networks trained using MSE $D = \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2$ in Eq. 3, and *factorized-msssim* networks use $D = 1 - (\text{MS-SSIM})$ in Eq. 3.

Conjectures / Observations

(6) As D_p is related to $p(x)$, \Rightarrow using $\mathcal{L}(\theta) = \int p(x) D_p(x, f(x)) dx$ implies considering $p(x)$ twice

(6.A) GOOD When your data is poor using D_p regularizes



(a) $\min_{u(\mathbf{x})} \|\mathbf{x} - \hat{\mathbf{x}}\|_2$



(b) $\min_{u(\mathbf{x})} \text{NLPD}(\mathbf{x}, \hat{\mathbf{x}})$



(c) $\min_{u(\mathbf{x})} 1 - \text{MS-SSIM}(\mathbf{x}, \hat{\mathbf{x}})$

Figure 6: **Checking observation 8.** Decoded image (compressed at 0.25bpp) encoded with networks trained using data from a uniform distribution and Euclidean vs perceptual losses.

Conjectures / Observations

(6) As \mathcal{D}_p is related to $p(x)$, \Rightarrow using $\mathcal{L}(\theta) = \int p(x) \mathcal{D}_p(x, f(x)) dx$ implies considering $p(x)$ twice

(6.A) GOOD When your data is poor using \mathcal{D}_p regularizes

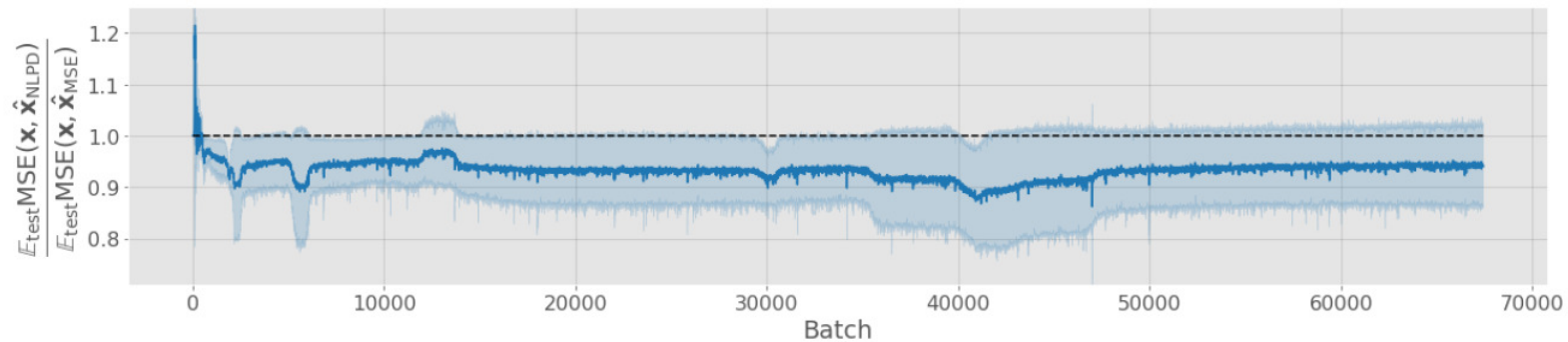


Figure 16: Gain in using NLPD over MSE as a loss function evaluated in terms of MSE loss on test set (Kodak dataset) using batch size of 1 and a small learning rate, fixing random seeds. \hat{x}_{NLPD} denotes the reconstruction of x with a network optimized for NLPD, and \hat{x}_{MSE} for a network optimized for MSE. The mean (solid line) and standard deviation (solid fill) was taken over 5 runs with different random seeds, i.e. different network initialization and training image ordering. The dashed line represents if the two networks had the same expected MSE on the test set.

Conjectures / Observations

(6) As \mathcal{D}_p is related to $p(x)$, \Rightarrow using $\mathcal{L}(\theta) = \int p(x) \mathcal{D}_p(x, f(x)) dx$ implies considering $p(x)$ twice

(6.B) BAD When training with \mathcal{D}_p performance may be poor in $p(x)$ \downarrow

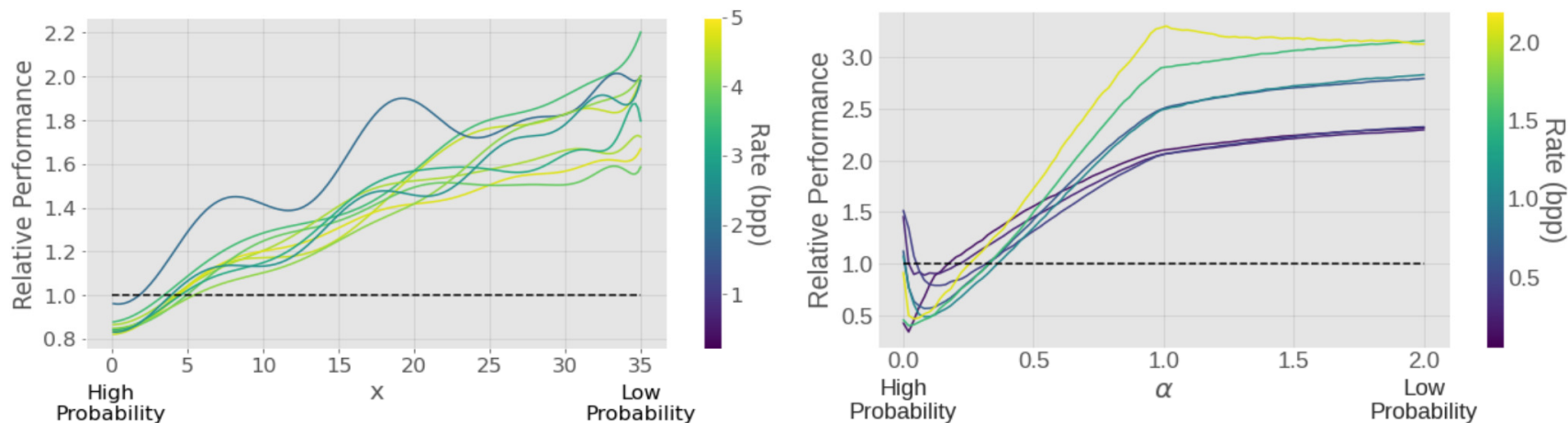


Figure 5: **Checking observation 7.** Relative performance of networks for samples along a line through the support of the respective distributions. Left: networks trained with $\mathcal{D} = p(\mathbf{x}) \cdot \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2$ in Eq. 3 divided by performance of networks trained with $\mathcal{D} = \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2$ on the 2D Student-t and evaluated using samples along x -axis (see Appendix B). Right: networks trained with $\mathcal{D} = 1 - (\text{MS-SSIM})$ in Eq. 3 divided by performance of networks trained with $\mathcal{D} = \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2$ for image coding using the Tensorflow Compression package.

CONCLUSIONS

① We considered the relations between

NATURAL
PERCEPTION

ARTIFICIAL
AUTOENCODERS

D_p Perceptual distance

D_r Reconstruction distance

D_i Inner domain distance

D_s Self-reconstr. distance

and PROBABILITY
 $P(\cdot)$

① Assuming the equalization property of the perceptual response, S , we proved $\frac{\partial D_p}{\partial x} = \frac{|S(x_1) - S(x_2)|_2}{|x_1 - x_2|_2} = P(x)^{\delta}$

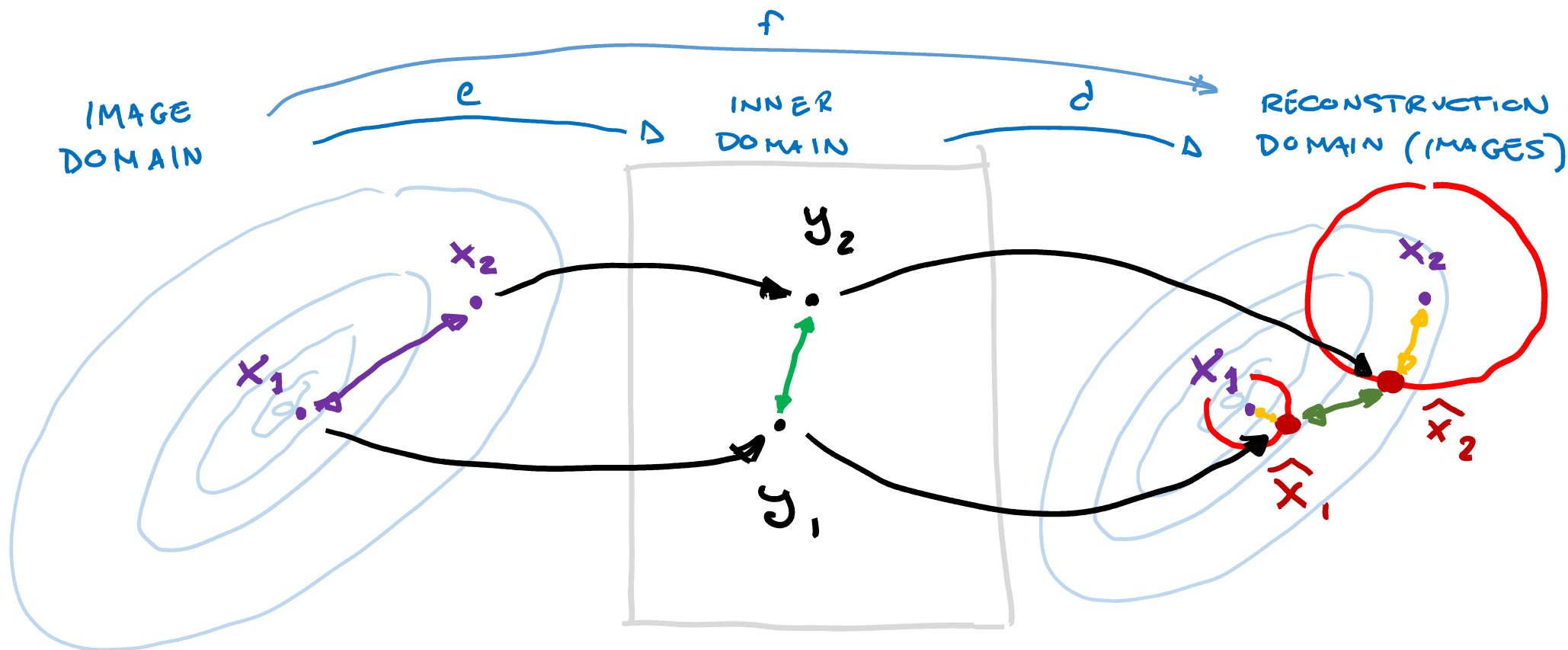
② Considering $\mathcal{L}(\theta) = \int P(x) D(x, f_{\theta}(x)) dx$ we conjecture $D_s \sim P(x)^{-1}$

③ Considering that the autoencoder response, $f(x)$, flattens $P(x)$ we conjecture $\frac{\partial D_r}{\partial x} = \frac{|f(x_1) - f(x_2)|_2}{|x_1 - x_2|_2} \sim P(x)$

④ Considering equalization in the inner domain we conjecture $\frac{\partial D_i}{\partial x} = \frac{|e(x_1) - e(x_2)|_2}{|x_1 - x_2|_2} \sim P(x)$

⑤ Assuming (1), (3) and (4) we conjecture D_r and D_i are correlated with D_p

⑥ Using D_p may imply taking into account $P(x)$ twice $\left\{ \begin{array}{l} 6.a \text{ Positive regularization effect of } D_p \\ 6.b \text{ Negative effect in some regions} \end{array} \right.$



EUCLIDEAN DISTANCE $D_e(x_1, x_2) = \|x_1 - x_2\|_2 = \text{RMSE}(x_1, x_2)$

INNER DISTANCE $D_i(x_1, x_2) = \|y_1 - y_2\|_2 = \|e(x_1) - e(x_2)\|_2$

RECONSTRUCT. DISTANCE $D_r(x_1, x_2) = \|\hat{x}_1 - \hat{x}_2\|_2 = \|f(x_1) - f(x_2)\|_2$

SELF RECONSTRUCT. DISTANCE $D_s(x, \hat{x}) = \|x - \hat{x}\|_2 = \|x - f(x)\|_2$