# Self Supervised Learning Methods for Imaging

## Part 2: Learning from noisy data

*Mike Davies, University of Edinburgh*

*Julián Tachella, CNRS, École Normale Supérieure de Lyon*

# Denoising problems

In this first part, we will focus on 'denoising' problems

$$\boldsymbol{y} = A\boldsymbol{x} + \boldsymbol{\epsilon}$$

where $A \in \mathbb{R}^{m \times n}$ is invertible (and thus $m \geq n$).

- We focus on $\boldsymbol{A} = \boldsymbol{I}$ for simplicity.

- All methods in this part can be extended to any invertible $\boldsymbol{A}$.

# Unsupervised Risk Estimators

**Supervised loss**

$$\mathcal{L}_{\text{sup}}(\boldsymbol{x}, \boldsymbol{y}, f) = ||\boldsymbol{x} - f(\boldsymbol{y})||^2 = ||\boldsymbol{y} - f(\boldsymbol{y})||^2 + 2f(\boldsymbol{y})^\top(\boldsymbol{y} - \boldsymbol{x}) + \text{const.}$$

Measurement consistency

key term to approximate!
$$= f(\boldsymbol{y})^\top \boldsymbol{\epsilon}$$

**Goal:** build a self-supervised loss $\mathcal{L}_{\text{self}}$ such that

$$\mathbb{E}_{\boldsymbol{y}} \, \mathcal{L}_{\text{self}}(\boldsymbol{y}, f) = \mathbb{E}_{\boldsymbol{x}, \boldsymbol{y}} \, \mathcal{L}_{\text{sup}}(\boldsymbol{x}, \boldsymbol{y}, f) + \text{const.}$$

# Noise2Noise

**Mallows $C_p$** [Mallows, 1973], **Noise2Noise** [Lehtinen, 2018]

- **Independent** pairs $y_a = x + \epsilon_a$ and $y_b = x + \epsilon_b$ with $\epsilon_a, \epsilon_b$ **independent**
- $\mathbb{E}_{\epsilon_b|x} \epsilon_b = 0$

$$\mathbb{E}_{y_b|x} f(y_a)^\top (y_b - x) = f(x + \epsilon_a) \underbrace{\mathbb{E} \, \epsilon_b}_{=0} = 0$$

$$\boxed{\mathcal{L}_{N2N}(y, f) = || \, y_b - f(y_a) ||^2}$$

- Also works for any noise distribution with $\mathbb{E}_{y_b|x} y_b = x$
- **Limitation:** observing independent copies is often impossible

# Noisier2Noise

**Recorrupted2Recorrupted** [Pang et al., 2021]**, Coupled Bootstrap** [Oliveira et al., 2022], **Noisier2Noise** [Moran et al., 2020].

**Proposition:** Let $\boldsymbol{y} \sim N(\boldsymbol{x}, I\sigma^2)$ and define

$$\boldsymbol{y}_a = \boldsymbol{y} + \alpha\boldsymbol{\omega}$$
$$\boldsymbol{y}_b = \boldsymbol{y} - \boldsymbol{\omega}/\alpha$$

where $\boldsymbol{\omega} \sim N(\boldsymbol{0}, I\sigma^2)$ and $\alpha \in \mathbb{R}$, then $\boldsymbol{y}_a$ and $\boldsymbol{y}_b$ are **independent** random variables (fixed $\boldsymbol{x}$).

$$\boxed{\mathcal{L}_{R2R}(\boldsymbol{y}, f) = \mathbb{E}_{\boldsymbol{\omega}}|| \boldsymbol{y}_b - f(\boldsymbol{y}_a) ||^2}$$

- Price to pay: SNR($\boldsymbol{y}_a$) < SNR($\boldsymbol{y}$)
- Trick can be extended to Poisson noise [Oliveira et al., 2023]
- At **test time**, $f^{\text{test}}(\boldsymbol{y}) = \frac{1}{N}\sum_i f(\boldsymbol{y} + \alpha\boldsymbol{\omega}_i)$ with $\boldsymbol{\omega}_i \sim \mathcal{N}(\boldsymbol{0}, I\sigma^2)$

# Stein's Unbiased Risk Estimator

- **Stein's lemma** [Stein 1974] **:** Let $\boldsymbol{y}|\boldsymbol{x} \sim \mathcal{N}(\boldsymbol{x}, I\sigma^2)$, $f$ be weakly differentiable, then

$$\mathbb{E}_{\boldsymbol{y}|\boldsymbol{x}} (\boldsymbol{y} - \boldsymbol{x})^\top f(\boldsymbol{y}) = \mathbb{E}_{\boldsymbol{y}|\boldsymbol{x}} \sigma^2 \sum_i \frac{\delta f_i}{\delta y_i}(\boldsymbol{y})$$

$$\mathcal{L}_{SURE}(\boldsymbol{y}, f) = ||\boldsymbol{y} - f(\boldsymbol{y})||^2 + 2\sigma^2 \sum_i \frac{\delta f_i}{\delta y_i}(\boldsymbol{y})$$

Measurement consistency      Degrees of freedom [Efron, 2004]

- **Hudson's lemma** [Hudson 1978] extends this result for the exponential family (eg. **Poisson Noise**)
- Beyond exponential family: **Poisson-Gaussian noise** [Le Montagner et al., 2014] [Raphan and Simoncelli, 2011]

# Stein's Unbiased Risk Estimator

**Monte Carlo SURE** [Efron 1975, Breiman 1992, Ramani et al., 2007]

SURE's divergence is generally approximated as

$$\sum_i \frac{\delta f_i}{\delta y_i}(\boldsymbol{y}) \approx \frac{\boldsymbol{\omega}^\top}{\alpha}\left(f(\boldsymbol{y}) - f(\boldsymbol{y} + \boldsymbol{\omega}\alpha)\right)$$

where $\alpha > 0$ small, $\boldsymbol{\omega} \sim \mathcal{N}(\boldsymbol{0}, I)$

- Noisier2Noise is equivalent to SURE when $\alpha \to 0$ [Oliveira, 2022].

# Stein's Unbiased Risk Estimator

The solution to SURE is **Tweedie's Formula**

$$\underset{f}{\arg\min} \; \mathbb{E}_{\boldsymbol{y}}|| \boldsymbol{y} - f(\boldsymbol{y})||^2 + 2\sigma^2 \sum_i \frac{\delta f_i}{\delta y_i}(\boldsymbol{y})$$

Integration by parts

$$\underset{f}{\arg\min} \; \mathbb{E}_{\boldsymbol{y}} || \boldsymbol{y} - f(\boldsymbol{y})||^2 - 2\sigma^2 \sum_i f(\boldsymbol{y}) \frac{\delta \log p_{\boldsymbol{y}}(\boldsymbol{y})}{\delta y_i}$$

Complete squares

$$\underset{f}{\arg\min} \; \mathbb{E}_{\boldsymbol{y}} || f(\boldsymbol{y}) - \boldsymbol{y} - \sigma^2 \nabla \log p_{\boldsymbol{y}}(\boldsymbol{y}) ||^2$$

$$\Longrightarrow \quad f(\boldsymbol{y}) = \boldsymbol{y} + \sigma^2 \nabla \log p_{\boldsymbol{y}}(\boldsymbol{y})$$

- **Noise2Score** [Kim and Ye, 2021] learns $\nabla \log p_{\boldsymbol{y}}(\boldsymbol{y})$ from noisy data + denoises with Tweedie.
- Key formula behind diffusion models, which can be trained self-supervised [Daras et al., 2024]

# Summary So Far

| | Train Eval | Test Eval | Single $y$ | MMSE optimal | Unknown noise |
|---|---|---|---|---|---|
| Noise2Noise | 1 | 1 | | ✅ | ✅ |
| Noisier2Noise | 1 | >1 | ✅ | ✅ | |
| SURE | 2 | 1 | ✅ | ✅ | |

If we have a single $y$ and don't know the noise distribution?

# Cross-Validation Methods

**Assumption:** $f_i$ does not depend on $y_i$, that is $\frac{\delta f_i}{\delta y_i} = 0$. Decomposable noise $p(\boldsymbol{y}|\boldsymbol{x}) = \prod p(y_i|x_i)$

$$\mathbb{E}_{\boldsymbol{y}|\boldsymbol{x}} \sum_{i=1}^{n} f_i(\boldsymbol{y})(y_i - x_i) = \sum_{i=1}^{n} \mathbb{E}_{\boldsymbol{y}_{-i}|\boldsymbol{x}} f_i(\boldsymbol{y}_{-i}) \overbrace{\mathbb{E}_{y_i|x_i} (y_i - x_i)}^{=0} = 0$$

$$\boxed{\mathcal{L}_{CV}(\boldsymbol{y}, f) = ||\, \boldsymbol{y} - f(\boldsymbol{y})||^2 \text{ subject to } \frac{\delta f_i}{\delta y_i}(\boldsymbol{y}) = 0 \ \forall i}$$

- SURE's perspective:

$$\mathcal{L}_{SURE}(\boldsymbol{y}, f) = ||\, \boldsymbol{y} - f(\boldsymbol{y})||^2 + 2\sigma^2 \sum_{i} \frac{\delta f_i}{\delta y_i}(\boldsymbol{y})$$

- These methods are not MMSE optimal
- How to remove dependence on $y_i$: training or architecture

# Measurement Splitting

**Cross-validation** [Efron, 2004]: random split $\boldsymbol{y} = \begin{bmatrix} \boldsymbol{y}_a \\ \boldsymbol{y}_b \end{bmatrix}$ at each iteration

$$\mathcal{L}_{N2V}(\boldsymbol{y}, f) = \mathbb{E}_{a,b} || \boldsymbol{y}_b - \text{diag } \boldsymbol{m}_b \, f(\boldsymbol{y}_a)||^2$$

where $\boldsymbol{m}_b \in \{0,1\}^n$ masks out the pixels in $\boldsymbol{y}_a$.

**Noise2Void** [Krull et al., 2019], **Noise2Self** [Batson, 2019]

- During training flip centre pixel
- Computes loss only on flipped pixels



**Neighbor2Neighbor** [Huang, 2023]

- Use different subsampling as input and target
- Assumes scale invariance

# Measurement Splitting

At **test time**, $f(\boldsymbol{y})$ is evaluated as

1. Test $f$ as trained (expensive)

$$f^{\text{test}}(\boldsymbol{y}) = \frac{1}{N}\Sigma_i M\, f(\boldsymbol{y}_{\boldsymbol{a}_i}) \text{ with } \boldsymbol{y}_{\boldsymbol{a}_i} \sim p(\boldsymbol{y}_{\boldsymbol{a}}|\boldsymbol{y}) \text{ and } M = \left(\Sigma_i^N \text{diag}\,(\boldsymbol{m}_{b,i})\right)^{-1}$$

2. Assume good generalization of $f$ (cheap)

- $f^{\text{test}}(\boldsymbol{y}) = f(\boldsymbol{y}_{\boldsymbol{a}})$ with $\boldsymbol{y}_{\boldsymbol{a}} \sim p(\boldsymbol{y}_{\boldsymbol{a}}|\boldsymbol{y})$

- $f^{\text{test}}(\boldsymbol{y}) = f(\boldsymbol{y})$

# Blind Spot Networks

**Blind spot networks** [Laine et al., 2019], [Lee et al., 2022]

- Convolutional architecture that doesn't 'see' centre pixel by construction

$$\mathcal{L}_{\text{BS}}(\boldsymbol{y}, f_{\text{BS}}) = || \boldsymbol{y} - f_{\text{BS}}(\boldsymbol{y})||^2$$

# Autoencoders

**Autoencoders**

**Assume**
- $f$ has a strong bottleneck

$$f_{\text{AE}} = d \cdot e$$



$$\mathbb{R}^n \quad e \quad \mathbb{R}^k \quad d \quad \mathbb{R}^n$$

$\sum_i \frac{\delta f_i}{\delta y_i}(\boldsymbol{y}) \approx O(k) \ll n$ is small

$$\boxed{\mathcal{L}_{AE}(\boldsymbol{y}, f) = || \boldsymbol{y} - f_{\text{AE}}(\boldsymbol{y})||^2}$$

- Noise distribution is 'high-dimensional' whereas signal distribution is 'low-dimensional'
- **Example:** linear ortho denoiser $f(\boldsymbol{y}) = M\boldsymbol{y}$, then $\sum_i \frac{\delta f_i}{\delta y_i}(\boldsymbol{y}) = \text{tr}\, M = k$

# Summary

| | Train Eval | Test Eval | Single $y$ | MMSE optimal | Unknown separable noise | Unknown coloured noise |
|---|---|---|---|---|---|---|
| Noise2Noise | 1 | 1 | | ✅ | ✅ | ✅ |
| Noisier2Noise | 1 | >1 | ✅ | ✅ | | |
| SURE | 2 | 1 | ✅ | ✅ | | |
| Noise2Void | 1 | 1 | ✅ | | ✅ | |
| Blind Spot | 1 | >1 | ✅ | | ✅ | |
| Autoencoders | 1 | 1 | ✅ | | ✅ | ✅ |

**No free lunch**: less assumptions about noise = less optimal estimator

# Beyond Denoising

For $A \neq I$, most estimators can be adapted to approximate

$$\mathbb{E}_{x,y} ||A^\dagger A(x - f(y))||^2$$

where $A^\dagger$ is the pseudoinverse of $A$.
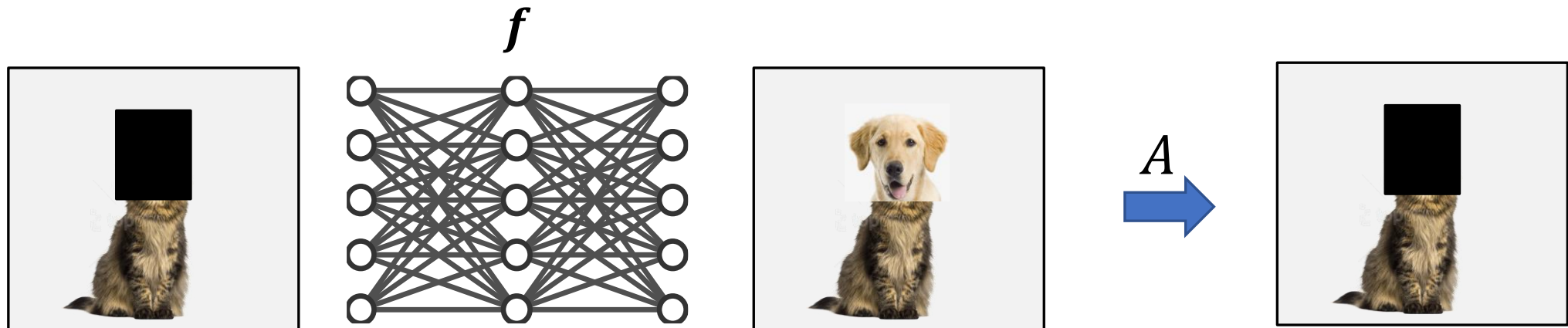
For example, **GSURE** [Eldar, 2008] writes for Gaussian noise

$$\mathcal{L}_{GSURE}(y, f) = ||A^\dagger y - A^\dagger A f(y)||^2 + 2\sigma^2 \sum_i \frac{\delta [A^\dagger A \cdot f]_i}{\delta y_i}(y)$$

# Incomplete Measurements?

1. If $A$ is invertible, we have $A^\dagger A = I$

2. If $A$ is not invertible, $\mathbb{E}_{x,y} ||A^\dagger A(x - f(y))||^2 \neq \mathbb{E}_{x,y} ||x - f(y)||^2$

In this case, the risk does not penalise $f(y)$ in the **nullspace** of $A$!

# References

The full reference list for this tutorial can be found here:

[https://tachella.github.io/projects/selfsuptutorial/](https://tachella.github.io/projects/selfsuptutorial/)