RED: <u>Regularization by Denoising</u> The Little Engine that Could

Michael Elad

May 18, 2017





Joint work with:





Yaniv Romano

Peyman Milanfar

This work was done during the summer 2016 in Google-Research Mountain-View, CA

Background and Main Objective



Image Denoising – Past Work

Searching Web-of-Science (May 4th, 2017) with

TOPIC = image and noise and (denoising or removal) Leads to ~4700 papers





Image Denoising – What's Out There?



Is Denoising Dead?

To a large extent, removal of additive noise from an image is a solved problem in image processing

This claim is based on several key observations:

- □ There are more than 20,000 papers studying this problem
- □ The proposed methods in the past decade are highly impressive
- □ Very different algorithms tend to lead to quite similar output quality
- □ Work investigating performance bounds confirm this [Chatterjee & Milanfar, 2010], [Levin & Nadler, 2011]

Bottom line: Improving image denoising algorithms seems to lead to diminishing returns

Noise Removal, Sharpening





Noise Removal, Sharpening





Compression Artifact Removal





Compression Artifact Removal





Some Visual Effects

Old man



Younger



More fun



Image Denoising – Can we do More?

SO

if improving image denoising algorithms is a dead-end avenue ... THFN

Lets seek ways to leverage existing denoising algorithms in order to solve OTHER (INVERSE) PROBLEMS

Inverse Problems ?

- Given the measurement y, recover its origin signal, x
- □ A statistical tie between x and y is given
- Examples: Image deblurring, inpainting, tomographic reconstruction, super-resolution, ...
- **Classic inverse problem:** $\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{v}$ where $\mathbf{v} \sim \mathbb{N}\left\{\mathbf{0}, \sigma^2 \mathbf{I}\right\}$
- The Bayesian approach: use the conditional probability Prob(x|y) in order to estimate x. How?
 - Maximize the conditional probability (MAP)
 - Find the expected x given y

"The Little Engine that Could" ?

As we are about to see, image denoisers (we refer to these as "engines") are capable of far more than just ... denoising, And thus this sub-title

Prior-Art 1:

Laplacian Regularization

Pseudo-Linear Denoising

Often times we can describe our denoiser as a pseudolinear Filter
Filter $\{x\} = W(x)x$

True for K-SVD, EPLL, NLM, BM3D and other algorithms, where the overall processing is divided into a non-linear stage of decisions, followed by a linear filtering

Typical properties of W:

- Row stochastic: <u>1</u> is an eigenvector, spectral radius is 1
- Fast decaying eigenvalues
- Symmetry or Non-negativity of **W**? Not always ...

The Image Laplacian

Given this form:

Filter
$$\{x\} = W(x)x$$



We may propose an image-adaptive (-)Laplacian:

Laplacian
$$\{x\} = x - W(x)x$$
 The "residual"
= $(I - W(x))x$
= $L(x)x$

Laplacians as Regularization Log-Likelihood Data Fidelity $\ell(\mathbf{x},\mathbf{y}) + \frac{\lambda}{2} \mathbf{x}^{\mathsf{T}} \mathbf{L} \mathbf{x}$ Regularization $\rho(\mathbf{x}) = \frac{1}{2} \mathbf{x}^{\mathsf{T}} \mathbf{L} \mathbf{x} = \frac{1}{2} \mathbf{x}^{\mathsf{T}} \left(\mathbf{x} - \mathbf{W} \mathbf{x} \right)$ **Filter Residual** Signal

This idea appeared is a series of papers in several variations:

[Elmoataz, Lezoray, & Bougleux 2008] [Szlam, Maggioni, & Coifman 2008]
[Peyre, Bougleux & Cohen 2011] [Milanfar 2013] [Kheradmand & Milanfar 2014]
[Liu, Zhai, Zhao, Zhai, & Gao 2014] [Haque, Pai, & Govindu 2014] [Romano & Elad 2015]

Laplacians as Regularization $\ell(x,y) + \frac{\lambda}{2}x^{T}(x - Wx)$

The problems with this line of work are that:

- The regularization term is hard to work with since L/W is a function of x. This is circumvented by cheating and assuming a fixed W per each iteration
- 2. If so, what is really the underlying energy that is being minimized?
- 3. When the denoiser cannot admit a pseudo-linear interpretation of W(x)x, this term is not possible to use

Prior-Art 2: The Plug-and-Play-Prior (P³) Scheme

The P³ Scheme

The idea of using a denoiser to solve general denoising problems was proposed under the name "Plug-and-Play-Priors" (P³) [Venkatakrishnan, Wohlberg & Bouman, 2013]

□ Main idea: Use ADMM to minimize the MAP energy

$$\hat{\mathbf{x}}_{\text{MAP}} = \min_{\mathbf{x}} \ell(\mathbf{x}, \mathbf{y}) + \frac{\lambda}{2} \rho(\mathbf{x})$$

$$\min_{\mathbf{x}, \mathbf{v}} \ell(\mathbf{x}, \mathbf{y}) + \frac{\lambda}{2} \rho(\mathbf{v}) \text{ s.t. } \mathbf{x} = \mathbf{v}$$



The above relies on a well-known concept in optimization called the augmented Lagrange algorithm, where u is the scaled Lagrange multiplier vector

The P³ Scheme $\min_{x,v} \ell(x,y) + \frac{\lambda}{2} \rho(v) + \frac{\beta}{2} \|x - v + u\|_{2}^{2}$

Minimize the above iteratively, in an alternating-direction fashion, w.r.t. the two unknowns (and update u=u-v+x):

1.
$$\min_{\mathbf{x}} \ell(\mathbf{x}, \mathbf{y}) + \frac{\beta}{2} \|\mathbf{x} - \mathbf{v} + \mathbf{u}\|_{2}^{2}$$
 : Simple inverse problem
2. $\min_{\mathbf{v}} \frac{\lambda}{2} \rho(\mathbf{v}) + \frac{\beta}{2} \|\mathbf{x} - \mathbf{v} + \mathbf{u}\|_{2}^{2}$



If the involved terms are convex, this algorithm is guaranteed to converge to the global optimum of the original function

Implicit Prior:

This is a key feature of P^3 – The idea that one can use ANY denoiser as a replacement for this stage, even if $\rho(v)$ is not known



If the involved terms are convex, this algorithm is guaranteed to converge to the global optimum of the original function

Implicit Prior:

This is a key feature of P^3 – The idea that one can use ANY denoiser as a replacement for this stage, even if $\rho(v)$ is not known

P³ Shortcomings

□ The P³ scheme is an excellent idea, but it has few troubling shortcomings:

- Parameter tuning is **TOUGH** when using a general denoiser
- This method is tightly tied to ADMM without an option for changing this scheme
- **CONVERGENCE** ? Unclear (steady-state at best)
- For an arbitrarily denoiser, no underlying & consistent
 COST FUNCTION

□ In this work we propose an alternative which is closely related to the above ideas (both Laplacian regularization and P³) which overcomes the mentioned problems: RED

RED: First Steps

Regularization by Denoising [RED]

We suggest:
$$\rho(\mathbf{x}) = \frac{1}{2}\mathbf{x}^{\mathsf{T}}(\mathbf{x} - \mathbf{W}\mathbf{x})$$

... for an arbitrary denoiser f(x)

1.
$$x = 0$$

 $\rho(x) = 0 \implies 2. x = f(x)$
3. Orthogonality

[Romano, Elad & Milanfar, 2016]

Regularization by Denoising [RED]

We suggest:
$$\rho(x) = \frac{1}{2}x^{T}(x-f(x))$$

... for an arbitrary denoiser f(x)

1.
$$x = 0$$

 $\rho(x) = 0 \implies 2. x = f(x)$
3. Orthogonality

[Romano, Elad & Milanfar, 2016]

Which f(x) to Use ?

Almost any algorithm you want may be used here, from the simplest Median (see later), all the way to the state-of-the-art CNN-like methods

We shall require f(x) to satisfy several properties as follows ...

Denoising Filter Property I $f(x):[0,1]^n \rightarrow [0,1]^n$

Differentiability:

- Some filters obey this requirement (NLM, Bilateral, Kernel Regression, TNRD)
- Others can be ε-modified to satisfy this (Median, K-SVD, BM3D, EPLL, CNN, ...)

Denoising Filter Property II

□ Local Homogeneity: for $|c-1| \le \varepsilon << 1$, we have that f(cx) = cf(x)



Are Denoisers Homogenous ?



Implication (1)

Directional Derivative:



Homogeneity

Looks Familiar ?

U We got the property

$$\nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x}) \cdot \mathbf{x} = \mathbf{f}(\mathbf{x})$$

This is much more general than f(x) = W(x)x and applies to any denoiser satisfying the above conditions

Implication (2)



Implication: Filter stability. Small additive perturbations of the input don't change the filter matrix
Denoising Filter Property III

Passivity via the spectral radius:

$$r\left\{\nabla_{x}f(x)\right\} = \max \left|\lambda\left(\nabla_{x}f(x)\right)\right| \leq 1$$
$$\implies \|x\| \geq r\left\{\nabla_{x}f(x)\right\} \cdot \|x\| \geq \|\nabla_{x}f(x)x\| = \|f(x)\|$$



Are Denoisers Passive ?

A direct inspection of this property implies computing the Jacobian of f(x) and evaluating its spectral radius r by the Power-Method:

$$\mathbf{h}_{k+1} = \frac{\nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x}) \cdot \mathbf{h}_{k}}{\left\| \nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x}) \cdot \mathbf{h}_{k} \right\|}$$

□ After enough iterations, $\mathbf{r} = \left| \mathbf{h}_{k+1}^{\mathsf{T}} \cdot \mathbf{h}_{k} \right|$ □ The problem is getting the Jacobian

Are Denoisers Passive ?

□ The alternative is based on the Taylor expansion

$$f(x+h) \cong f(x) + \nabla_x f(x) \cdot h \Longrightarrow \nabla_x f(x) \cdot h = f(x+h) - f(x)$$

The Power Method becomes: $h_{k+1} = \frac{\nabla_x f(x) \cdot h_k}{\left\| \nabla_x f(x) \cdot h_k \right\|} \longrightarrow h_{k+1} = \frac{f(x+h_k) - f(x)}{\left\| f(x+h_k) - f(x) \right\|}$

This leads to a value smaller than 1 for K-SVD, BM3D, NLM, EPLL, and the TNRD

Summary of Properties

□ The 3 properties that f(x) should follow:



RED: Advancing

Regularization by Denoising (RED) $\rho(x) = \frac{1}{2}x^{T}(x - f(x))^{*}$

Surprisingly, this expression is differentiable:

*

$$\nabla \rho(\mathbf{x}) = \mathbf{x} - \frac{1}{2} \nabla \left\{ \mathbf{x}^{\mathsf{T}} f(\mathbf{x}) \right\}$$
$$= \mathbf{x} - \frac{1}{2} \left(f(\mathbf{x}) + \nabla f(\mathbf{x}) \mathbf{x} \right) = \mathbf{x} - f(\mathbf{x}) \text{ the residual}$$
$$\sum_{\mathsf{Why not } \rho(\mathsf{x}) = \frac{1}{2} \| \mathbf{x} - f(\mathbf{x}) \|_{2}^{2}} \qquad \nabla_{\mathsf{x}} f(\mathbf{x}) \cdot \mathbf{x} = f(\mathsf{x}) \text{ and Homogeneity}$$

Regularization by Denoising (RED) $\rho(\mathbf{x}) = \frac{1}{2} \mathbf{x}^{\mathsf{T}} \left(\mathbf{x} - \mathbf{f}(\mathbf{x}) \right)$ $\nabla \rho(\mathbf{x}) = \mathbf{x} - \mathbf{f}(\mathbf{x})$ $\nabla \left\{ \nabla \rho(\mathbf{x}) \right\} = \mathbf{I} - \nabla_{\mathbf{x}} \mathbf{f}(\mathbf{x}) \succeq \mathbf{0}$ Relying on the differentiability Passivity guarantees positive $r{\nabla_{x}f(x)} \leq 1$ definiteness of the Hessian and hence convexity

RED for Inverse Problems

$$\begin{array}{l} \underset{x}{\text{min }} \ell(x,y) + \frac{\lambda}{2} x^{\mathsf{T}} (x - f(x)) \\ \text{Log-Likelihood} \\ \text{Data Fidelity} \end{array}$$

Our regularization term is convex and thus the whole expression is convex if the Log-likelihood term is convex as well

RED for Linear Inverse Problems

$$\min_{x} \frac{1}{2} \|\mathbf{H}\mathbf{x} - \mathbf{y}\|_{2}^{2} + \frac{\lambda}{2} \mathbf{x}^{\mathsf{T}} (\mathbf{x} - \mathbf{f}(\mathbf{x}))$$

L₂-based Data Fidelity Regularization

This energy-function is convex

Any reasonable optimization algorithm will get to the global minimum if applied correctly

Numerical Approach: Gradient Descent

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{H}\mathbf{x} - \mathbf{y}\|_{2}^{2} + \frac{\lambda}{2} \mathbf{x}^{\mathsf{T}} (\mathbf{x} - \mathbf{f}(\mathbf{x}))$$
$$= \left(\mathbf{I} - \mu \left[\mathbf{H}^{\mathsf{T}}\mathbf{H} + \lambda \mathbf{I}\right]\right) \mathbf{x}_{k} + \mu \left(\mathbf{H}^{\mathsf{T}}\mathbf{y} + \lambda \mathbf{f}(\mathbf{x}_{k})\right)$$

Guaranteed to converge for $0 < \mu < 1$, when **H** is stable

Drawback: Every iteration applies f(x) but mildly "inverting" the blur effect

Numerical Approach: Gradient Descent



Numerical Approach II: ADMM $\min_{x} \frac{1}{2} \|Hx - y\|_{2}^{2} + \frac{\lambda}{2} x^{T} (x - f(x))$

□ First: apply variable splitting by setting x = v in the second term:

$$\min_{x,v} \frac{1}{2} \| \mathbf{H} x - y \|_{2}^{2} + \frac{\lambda}{2} v^{\mathsf{T}} (v - f(v)) \quad \text{s.t. } x = v$$

□ Second: Use Augmented Lagrangian:

$$\min_{x,v} \frac{1}{2} \|\mathbf{H}x - y\|_{2}^{2} + \frac{\lambda}{2} v^{T} (v - f(v)) + \frac{\beta}{2} \|x - v + u\|_{2}^{2}$$

Numerical Approach II: ADMM

□ Third: Optimize w.r.t. x and v alternatingly

• Update x:
$$\min_{x} \frac{1}{2} \|Hx - y\|_{2}^{2} + \frac{\beta}{2} \|x - v + u\|_{2}^{2} \quad \text{Solve a simple} \\ \text{Inear system} \\ \text{Update } v: \min_{v} \frac{\lambda}{2} v^{\mathsf{T}} (v - f(v)) + \frac{\beta}{2} \|x - v + u\|_{2}^{2} \\ \implies \lambda (v - f(v)) - \beta (x - v + u) = 0 \quad \text{Solve by fixed point iteration} \\ \lambda (v_{k+1} - f(v_{k})) + \beta (x - v_{k+1} + u) = 0 \implies v_{k+1} = \frac{\lambda f(v_{k}) + \beta (x + u)}{\lambda + \beta} \\ \end{cases}$$

The P³ differs from this ADMM in this stage, offering to compute $v_{k+1} = f(x+u)$

Numerical Approach II: ADMM

□ Third: Optimize w.r.t. x and v alternatingly

Update x:
$$\min_{x} \frac{1}{2} \| Hx - y \|_{2}^{2} + \frac{\beta}{2} \| x - v + u \|_{2}^{2}$$
Solve a simple linear system
Update v: $\min_{v} \frac{\lambda}{2} v^{\mathsf{T}} (v - f(v)) + \frac{\beta}{2} \| x - v + u \|_{2}^{2}$
 $\Longrightarrow \lambda (v - f(v)) - \beta (x - v + u) = 0$
Solve by fixed point iteration
 $\lambda (v_{k+1} - f(v_{k})) + \beta (x - v_{k+1} + u) = 0$
 $\searrow v_{k+1} = \frac{\lambda f(v_{k+1} - f(v_{k}))}{\zeta u^{2} r^{2} r^{2}$

The P³ differs from this ADMM in this stage, offering to compute $v_{k+1} = f(x+u)$

Numerical Approach: Fixed Point



Numerical Approach III: Fixed Point





$$z_{k+1} = f(Mz_k + b)$$

- ❑ While CNN use a trivial and weak nonlinearity f(●), we propose a very aggressive and image-aware denoiser
- Our scheme is guaranteed to minimize a clear and relevant objective function



$$z_{k+1} = f(Mz_k + b)$$

- ❑ While CNN use a trivial and weak nonlinearity f(●), we propose a very aggressive and image-aware denoiser
- Our scheme is guaranteed to minimize a clear and relevant objective function

RED Underlying Model

RED assumes the following Prior

$$Prob(x) = C \cdot exp\{-c \cdot x^{T}(x - f(x))\}$$

Theorem: With respect to the above prior, f(x) is the MMSE estimator

Conjecture: With respect to this prior, minimizing $\min_{x} \frac{1}{2} \|\mathbf{H}\mathbf{x} - \mathbf{y}\|_{2}^{2} + \frac{\lambda}{2} \mathbf{x}^{\mathsf{T}} (\mathbf{x} - \mathbf{f}(\mathbf{x}))$ gives the optimal MMSE solution

So, again, Which f(x) to Use ?

Almost any algorithm you want may be used here, from the simplest Median (see later), all the way to the state-of-the-art CNN-like methods

Comment: Our approach has one hidden parameter – the level of the noise (σ) the denoiser targets. We simply fix this parameter for now. But more work is required to investigate its effect

RED and P³ Equivalence

RED vs. P^3 ?

- ADMM can be used for both P³ and RED, leading to two similar (yet different) algorithms
- These algorithms differ in the update stage of v (assign y=x+u):
 - P³: v=f(y)
 - RED: v is the solution of $\lambda(v-f(v))-\beta(y-v)=0$
- **Question**: Under which conditions on f(x) (and λ) would the two be the same?

RED for Denoising?

- Assume f(x)=Wx in these equations, in order to simplify them:
 - P³: v=f(y)=Wy
 - RED: λ(v-f(v))-β(y-v)=λ(v-Wv)-β(y-v)=0

Thus, for getting the same outcome we require $\lambda (\mathbf{v} - \mathbf{W}\mathbf{v}) - \beta (\mathbf{y} - \mathbf{v}) = 0 \quad \& \quad \mathbf{v} = \mathbf{W}\mathbf{y}$ $\implies \lambda (\mathbf{W}\mathbf{y} - \mathbf{W}\mathbf{W}\mathbf{y}) - \beta (\mathbf{y} - \mathbf{W}\mathbf{y}) = 0$

RED for Denoising?

□ We got the equation

$$\lambda (\mathbf{W}\mathbf{y} - \mathbf{W}\mathbf{W}\mathbf{y}) - \beta (\mathbf{y} - \mathbf{W}\mathbf{y}) = \mathbf{0}$$

$$\begin{bmatrix} \beta \mathbf{I} - (\beta + \lambda)\mathbf{W} + \lambda \mathbf{W}^{2} \end{bmatrix} \mathbf{y} = (\beta \mathbf{I} - \lambda \mathbf{W})(\mathbf{I} - \mathbf{W})\mathbf{y} = \mathbf{0}$$

Answer: The two algorithms coincide if the eigenvalues of W are either 1 or β/λ, i.e. only when the denoiser is highly degenerated

Denoising via RED

RED for Denoising?

$$\min_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_{2}^{2} + \frac{\lambda}{2} \mathbf{x}^{\mathsf{T}} (\mathbf{x} - \mathbf{f}(\mathbf{x}))$$

$$\sum_{\mathbf{x} - \mathbf{y} + \lambda} (\mathbf{x} - \mathbf{y} + \lambda (\mathbf{x} - \mathbf{f}(\mathbf{x})) = 0$$

$$\int_{\mathbf{f}(\mathbf{x}) = \mathbf{W}\mathbf{x}} (\mathbf{x} - \mathbf{y} + \lambda (\mathbf{x} - \mathbf{w}))^{-1} \mathbf{y}$$

This is a "sharper" version of the original denoiser, known to modify the trading between bias and variance [Elmoataz et al. '08],[Bougleux et al. '09]

RED for Denoising?

Question: Could we claim this just gives Wy? (i.e. RED is equivalent to the denoiser it is built upon?)

Answer: Lets plug x=Wy into the equation

$$\mathbf{x} - \mathbf{y} + \lambda (\mathbf{x} - \mathbf{W}\mathbf{x}) = \mathbf{0} \Rightarrow \mathbf{W}\mathbf{y} - \mathbf{y} + \lambda (\mathbf{W}\mathbf{y} - \mathbf{W}\mathbf{W}\mathbf{y}) = \mathbf{0}$$
$$\begin{bmatrix} \mathbf{I} - (\mathbf{1} + \lambda)\mathbf{W} + \lambda \mathbf{W}^2 \end{bmatrix} \mathbf{y} = \mathbf{0}$$

The two filters are equivalent if the eigenvalues of **W** are either 1 or $1/\lambda$

Surprisingly: This condition is similar (equivalent for β =1) to the one leading to equivalence of RED and P³

RED in Practice

Examples: Deblurring



(d) NCSR 30.03dB (e) P^3 -TNRD 30.36dB (f) RED: ADMM-TNRD 30.40dB

Uniform 9×9 kernel and WAGN with $\sigma^2=2$

Examples: Convergence Comparison



Examples: Deblurring



Uniform 9×9 kernel and WAGN with $\sigma^2=2$



(d) NCSR 28.39dB

(e) P^3 -TNRD 28.43dB

(f) RED: FP-TNRD 28.82dB

Examples: 3x Super-Resolution



⁽a) Bicubic 20.68dB





Degradation:

- A Gaussian 7×7
 blur with width
 1.6
- A 3:1 downsampling and
 WAGN with σ=5

(d) Ours: SD-TNRD 27.39dB

(c) P³-TNRD 26.61dB

Examples: 3x Super-Resolution



(a) Bicubic 20.44dB



(c) P^3 -TNRD 23.25dB



(b) NCSR 22.97dB



(d) Ours: ADMM-TNRD 23.28dB

Degradation:

- A Gaussian 7×7
 blur with width
 1.6
- A 3:1 downsampling and
- WAGN with σ =5

Sensitivity to Parameters



(a) Comparison between RED and PPP



(b) P^3 : Sensitivity to change in β_0 and α

Sensitivity to Parameters



Conclusions
What have we Seen Today ?

- RED a method to take a denoiser and use it sequentially for solving inverse problems
- Main benefits: Clear objective being minimized, Convexity, flexibility to use almost any denoiser and any optimization scheme
- One could refer to RED as a way to substantiate earlier methods (Laplacian-Regularization and the P³) and fix them
- Challenges: Trainable version? Compression? MMSE conjecture?

Relevant Reading

- 1. "Regularization by Denoising", Y. Romano, M. Elad, and P. Milanfar, To appear, SIAM J. on Imaging Science.
- 2. "A Tour of Modern Image Filtering", P. Milanfar, IEEE Signal Processing Magazine, no. 30, pp. 106–128, Jan. 2013
- 3. "A General Framework for Regularized, Similarity-based Image Restoration", A. Kheradmand, and P. Milanfar, IEEE Trans on Image Processing, vol. 23, no. 12, Dec. 2014
- 4. "Boosting of Image Denoising Algorithms", Y. Romano, M. Elad, SIAM J. on Image Science, vol. 8, no. 2, 2015
- 5. "How to SAIF-ly Improve Denoising Performance", H. Talebi , X. Zhu, P. Milanfar, IEEE Trans. On Image Proc., vol. 22, no. 4, 2013
- 6. "Plug-and-Play Priors for Model-Based Reconstruction", S.V. Venkatakakrishnan, C.A. Bouman, B. Wohlberg, GlobalSIP, 2013
- 7. "BM3D-AMP: A New Image Recovery Algorithm Based on BM3D Denoising", C.A. Metzler, A. Maleki, and R.G. Baraniuk, ICIP 2015
- 8. "Is Denoising Dead?", P. Chatterjee, P. Milanfar, IEEE Trans. On Image Proc., vol. 19, no. 4, 2010
- 10. "Symmetrizing Smoothing Filters", P. Milanfar, SIAM Journal on Imaging Sciences, Vol. 6, No. 1, pp. 263–284
- 11. "What Regularized Auto-Encoders Learn from The Data-Generating Distribution", G. Alain, and Y. Bengio, JMLR, vol.15, 2014,
- 12. "Representation Learning: A Review and New Perspectives", Y. Bengio, A. Courville, and P. Vincent, IEEE Trans. on Pattern Analysis and Machine Intelligence, vol. 35, no. 8, Aug. 2013

Thank You

"That's all Folks!